

Random Sampling

Copyright Dick Startz

1

Random samples

A data point is a number (or a color or something). It might be a vector.

A data point is a *realization* of a draw from a *population*. We often assume a hypothetical population which is infinitely countable.

Copyright Dick Startz

2

Random samples

Random sample

The vector of random variables X_1, \dots, X_n is called a *random sample* of size n from the population $f(x)$ if X_1, \dots, X_n are mutually independent random variables and the marginal pdf of each X_i is the same function $f(x)$.

Alternatively, X_1, \dots, X_n are called *independent and identically distributed* random variables with pdf $f(x)$. This is commonly abbreviated *iid*.

Copyright Dick Startz

3

Joint density

$$f(x_1, \dots, x_n) = f(x_1) \times \dots \times f(x_n) = \prod_{i=1}^n f(x_i)$$

Not everything fits:

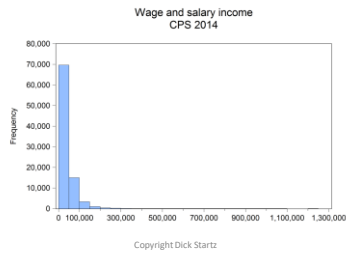
- The draws are not independent.
- The draws are not identical.
- Stratified sample.

Copyright Dick Startz

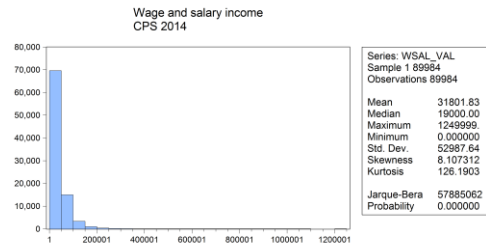
4

Histogram

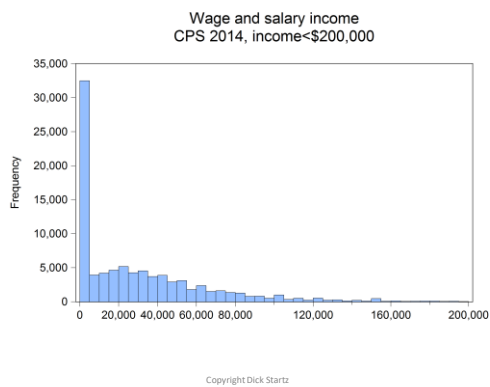
- Bar chart where each *bin* gives the count of observations within some range.



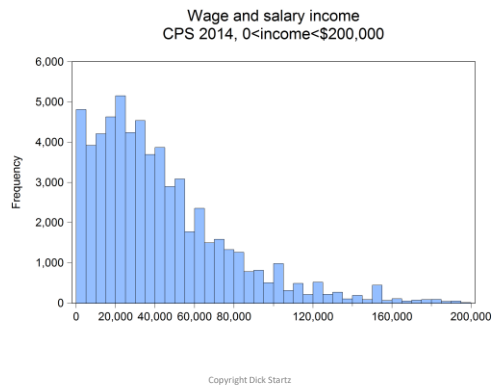
5



6



7



8

Histogram->pdf

Suppose we have B bins with counts c_i and widths w_i . Then for a histogram to represent a pdf we need a normalization constant k to make the cdf end up at 1.

$$1 = \sum_{b=1}^B k c_i w_i$$

$$k = \frac{1}{\sum_{b=1}^B c_i w_i}$$

New height is $k c_i$.

In the case of equal width bins this gives us

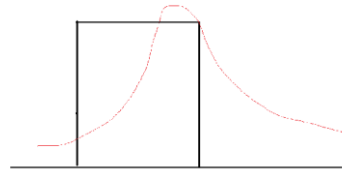
$$k = \frac{1}{wn}$$

Copyright Dick Startz

9

Kernel smoothing

- Wide bin leads to bias



Copyright Dick Startz

10

Narrow bin

A crude estimate of the density function around bin j is

$$\frac{n_j}{nh}$$

Where h is the width of the bin. n_j is distributed binomial with variance

$$\text{var}(n_j) = np(1-p)$$

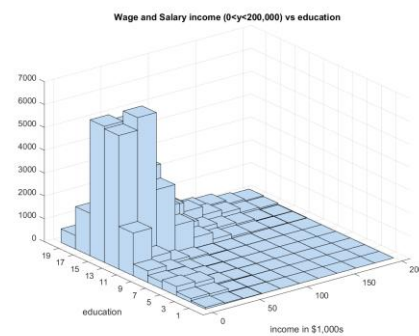
Where $p \propto f(\cdot)h$

So the variance of the fraction is

$$\frac{np(1-p)}{(nh)^2} \propto \frac{nf(\cdot)h(1-f(\cdot)h)}{n^2h^2}$$

Copyright Dick Startz

11



Copyright Dick Startz

12

Finite population

- Two balls in urn: one red, one black

Marginal distribution of a draw is $f(\text{black}) = 1/2$, but draws are not independent so two draws is not a “random sample.”

Copyright Dick Startz

13

Statistic

Let X_1, \dots, X_n be a random sample of size n from a population and let $T(x_1, \dots, x_n)$ be a real-valued or vector-valued function whose domain includes the sample space of (X_1, \dots, X_n) . Then the random variable or random vector $Y = T(X_1, \dots, X_n)$ is called a *statistic*. The probability distribution of a statistic Y is called the *sampling distribution* of Y .

Copyright Dick Startz

14

Examples of statistics

- Height of a bar in a histogram.
- Sample mean
- Sample variance
- Fraction of times a test is rejected in a Monte Carlo experiment.

Copyright Dick Startz

15

- Try the following problem. Draw m iid standard normals, $x_i, i = 1, \dots, m$. Compute the statistic

$$T = \frac{1}{m} \sum_{i=1}^m I(x_i < \Phi^{-1}(.025) \cup x_i > \Phi^{-1}(.975))$$

- Now repeat this experiment n times and show the distribution of T . Do this for $m = 100$ and $n = 1000$. Show what you get empirically as well as what the theoretical answer should be. (Remember that the test statistic is essentially an average of Bernoulli trials.)

Copyright Dick Startz

16

Sample mean and variance

The sample mean is

$$\bar{x} = \frac{x_1 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

And the sample variance is

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

The sample standard deviation is $s = \sqrt{s^2}$.

Copyright Dick Startz

19

Theorem 5.2.4: Let x_1, \dots, x_n be any numbers and $\bar{x} = (x_1 + \dots + x_n)/n$. Then

$$\begin{aligned} a. \quad \min_a \sum (x_i - a)^2 &= \sum (x_i - \bar{x})^2 \\ b. \quad (n-1)s^2 &\equiv \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \end{aligned}$$

Proof of part (a). Add and subtract \bar{x} and expand

$$\begin{aligned} \sum (x_i - a)^2 &= \sum ((x_i - \bar{x}) + (\bar{x} - a))^2 \\ &= \sum (x_i - \bar{x})^2 + 2 \sum (x_i - \bar{x})(\bar{x} - a) + \sum (\bar{x} - a)^2 \end{aligned}$$

The middle term equals zero because $(\bar{x} - a)$ is a constant and $\sum (x_i - \bar{x}) = 0$. The first term is positive. The last term is minimized at zero.

Copyright Dick Startz

20

Theorem 5.2.4: Let x_1, \dots, x_n be any numbers and $\bar{x} = (x_1 + \dots + x_n)/n$. Then

$$\begin{aligned} a. \quad \min_a \sum (x_i - a)^2 &= \sum (x_i - \bar{x})^2 \\ b. \quad (n-1)s^2 &\equiv \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \end{aligned}$$

The proof of (b) is

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i \bar{x} + \sum_{i=1}^n \bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - 2\bar{x}(n\bar{x}) + n\bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - n\bar{x}^2 \end{aligned}$$

Copyright Dick Startz

21

Unbiased

$\hat{\theta}$ is an unbiased estimator of θ if $E(\hat{\theta}) = \theta$.

Copyright Dick Startz

22

Theorem 5.2.6: Let x_1, \dots, x_n be a random sample from a population with mean μ and variance $\sigma^2 < \infty$.

- a. $E(\bar{x}) = \mu$, \bar{x} is *unbiased*.
- b. $\text{var}(\bar{x}) = \frac{\sigma^2}{n}$
- c. $E(s^2) = \sigma^2$, s^2 is *unbiased*.

$$E(\bar{x}) = \mu, \bar{x} \text{ is unbiased}$$

$$\begin{aligned} E(\bar{x}) &= E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \sum_{i=1}^n \mu \\ &= \mu \end{aligned}$$

Copyright Dick Startz

23

Copyright Dick Startz

24

$$\text{var}(\bar{x}) = \frac{\sigma^2}{n}$$

$$\begin{aligned} \text{var}(\bar{x}) &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} \text{var}\left(\sum_{i=1}^n x_i\right) \\ &= \frac{1}{n^2} \left(\sum_{i=1}^n \text{var } x_i\right) = \frac{1}{n^2} (n\sigma^2) = \frac{\sigma^2}{n} \end{aligned}$$

Copyright Dick Startz

25

$$E(s^2) = \sigma^2, s^2 \text{ is unbiased.}$$

$$\begin{aligned} E(s^2) &= E\left(\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right) \\ &= E\left(\frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - n\bar{x}^2\right]\right) \\ &= \frac{1}{n-1} \left(n(\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right)\right) \\ &= \frac{1}{n-1} (n\sigma^2 - \sigma^2) = \sigma^2 \end{aligned}$$

Copyright Dick Startz

26

Biased estimators

$$\tilde{s}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} s^2$$

$$s = (s^2)^{\frac{1}{2}}$$

is a concave function, so by Jensen's inequality

$$E(s) < [E(s^2)]^{\frac{1}{2}} = (\sigma^2)^{\frac{1}{2}} = \sigma$$

Copyright Dick Startz

27

Example of required sample size calculation

- Without treatment, outcome has mean zero.
- Wish to be able to detect effect of size 2, i.e. reject null of no effect by standard t -test.

How large a sample do you need?

Pilot study with no treatment shows outcome $\sim N(0, 120)$.

Copyright Dick Startz

28

Treatment data

$$\hat{\mu} \sim N\left(\mu, \frac{120}{n}\right)$$

Test of null:

$$z = \frac{\hat{\mu}}{\sqrt{120/n}}$$

Reject if $|z| > 1.96$.

Copyright Dick Startz

29

Probability of rejection

What is the probability that $|z| > 1.96$ if $\mu = 2$?

$$F_z(-1.96) + (1 - F_z(1.96))$$

$$z \sim N\left(\frac{\mu}{\sqrt{120/n}}, 1\right)$$

The cdf of z is

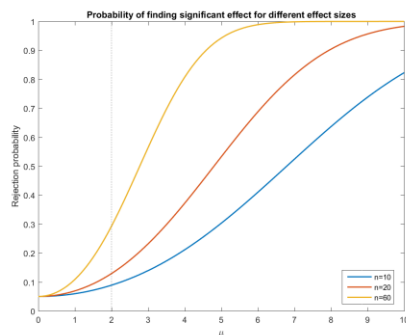
$$F_z(\cdot) = \Phi\left(z - \frac{\mu}{\sqrt{120/n}}\right)$$

Probability of rejection is

$$\Phi\left(-1.96 - \frac{\mu}{\sqrt{120/n}}\right) + \left(1 - \Phi\left(1.96 - \frac{\mu}{\sqrt{120/n}}\right)\right)$$

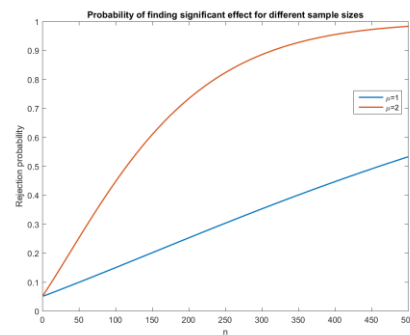
Copyright Dick Startz

30



Copyright Dick Startz

31



Copyright Dick Startz

32

My daughter is going to run a survey in Nigeria. Enumerators are expensive, so each marginal response costs $c = 2500$ naira. A field experiment is planned. A baseline survey was run before any treatment, which indicated that the standard deviation of the variable of interest is $\sigma = 4$ and that the responses are approximately normally distributed. My daughter believes the true effect size is $\mu = 2$. Finding a positive estimated effect size will result in a publication which will have a NPV for her career equal to \$50,000. (A reminder that sample means are roughly normally distributed.) How many surveys should my risk-neutral daughter buy?

Copyright Dick Startz

33