

Session 1 Introduction

What is **statistics**?

How to learn about a **population** given a **sample** of observations from that population.

- US non-institutional working civilians aged 16 or above in 2018. Interested in average hourly earnings in this population?
- All Singapore households in 2020. Interested in no. of dogs per SG household on average?

Math Review: Summation Notation

Given a set of numbers $\{x_i\}_{i=1}^n = \{x_1, x_2, \dots, x_n\}$, define

$$\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$$

Two Rules:

- $\sum_{i=1}^n (a_i + b_i) = \sum_{i=1}^n a_i + \sum_{i=1}^n b_i$
- $\sum_{i=1}^n ca_i = c \sum_{i=1}^n a_i$ where c is some constant value

Summation Notation

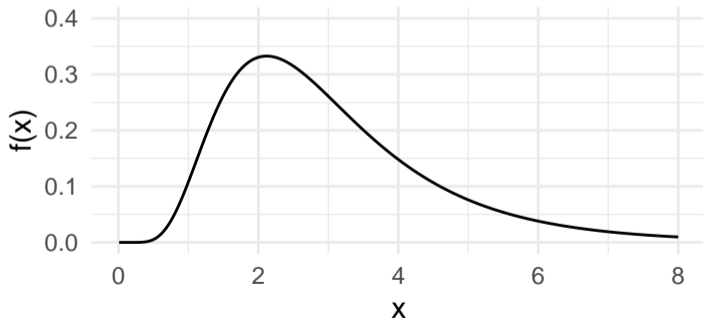
Proof of first equality

$$\begin{aligned}\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^n (x_i - \bar{x})y_i - \sum_{i=1}^n (x_i - \bar{x})\bar{y} \\ &= \sum_{i=1}^n (x_i - \bar{x})y_i - \bar{y} \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_{=0} \\ &= \sum_{i=1}^n (x_i - \bar{x})y_i\end{aligned}$$

Probability Review

$$X \sim \text{Log-normal}(\mu, \sigma^2)$$

$$X \sim \text{Log-normal}(\mu, \sigma^2) \iff \ln X \sim \text{Normal}(\mu, \sigma^2)$$



Probability Review: Expected Values

Mean or expected value of a random variable X is defined as

$$E(X) = \begin{cases} \sum_x x f_X(x) = \sum_x x \Pr(X = x) & \text{if } X \text{ is discrete, and} \\ \int_{-\infty}^{+\infty} x f_X(x) dx & \text{if } X \text{ is continuous.} \end{cases}$$

Measures “center” or “location” of distribution

- If $X \sim \text{Bernoulli}(p)$, then $E(X) = 1 \cdot p + 0 \cdot (1 - p) = p$.
- If $X \sim \text{Poisson}(\lambda)$, then

$$E(X) = \sum_{x=0}^{\infty} x \Pr(X = x) = \sum_{x=0}^{\infty} \frac{x e^{-\lambda} \lambda^x}{x!} = \lambda.$$

Probability Review: Expected Values

If X is a random variable, then $g(X)$ is also a random variable, with expectation:

$$E(g(X)) = \begin{cases} \sum_X g(x) f_X(x) = \sum_X g(x) \Pr(X = x) & \text{if } X \text{ is discrete, and} \\ \int_{-\infty}^{+\infty} g(x) f_X(x) dx & \text{if } X \text{ is continuous.} \end{cases}$$

For example,

- if $X \sim \text{Bernoulli}(p)$, we have $E(X^2) = 1^2p + 0^2(1 - p) = p$
- if $X \sim \text{Poisson}(\lambda)$, then $E(X^2) = \lambda + \lambda^2$
- if $X \sim \text{Normal}(\mu, \sigma^2)$, then $E(X^2) = \sigma^2 + \mu^2$.

Statistical Model

Often we **do not** need to specify distribution fully

We can assume

X_i iid such that $E(X_i) = \mu$ and $Var(X_i) = \sigma^2 < \infty$ for all $i = 1, \dots, n$.

NB: μ and σ^2 here is used as “generic notation” for mean and variance. E.g., if X_i is bernoulli, then $\mu = p$, $\sigma^2 = p(1 - p)$. If X_i is Poisson, then $\mu = \sigma^2 = \lambda$, etc.

Very general model! Assumes only that:

- sample is a random sample
- population is well-represented by *some* distribution with a finite mean and a finite variance

Bias

One commonly used criterion is **unbiasedness**: $E(\hat{\theta}) = \theta$

Sample mean is unbiased for true mean (under our stated conditions):

Proof:
$$E(\bar{X}) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} n\mu = \mu$$

- You will not *systematically* over- or under-estimate the population mean.
- (Thought experiment) If, say, 200 people went to the population and each obtained a random sample of n individuals and calculated the sample mean. Each would obtain a different sample mean, but their sample means will be nicely centered around the true (unknown) population mean.

Standard Error of Sample Mean

We should always try to get some idea of the potential size of estimation error

Consider variance of sample mean

$$Var(\bar{X}) = Var\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n Var(X_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}$$

- Not operational since σ^2 is unknown
- But already informative – variance goes down with increasing sample size

Standard Error of Sample Mean

To get numerical estimate of potential estimation error size:

- estimate σ^2 (let's call it $\widehat{\sigma}^2$ for now)

- Use square root of $\widehat{Var}(\bar{X}) = \frac{\widehat{\sigma}^2}{n}$, i.e., $\text{s.e.}(\bar{X}) = \sqrt{\frac{\widehat{\sigma}^2}{n}}$

Since $Var(X) = \sigma^2 = E((X - E(X))^2) = E(X^2) - E(X)^2$, it seems reasonable to use

$$\widetilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2$$

Unfortunately, this is a (downward) biased estimator for σ^2

Standard Error of Sample Mean

Fortunately, in this case, there is an obvious unbiased estimator. Let

$$\widehat{\sigma^2} = \frac{n}{n-1} \widetilde{\sigma^2} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad (\text{sample variance})$$

Then

$$E(\widehat{\sigma^2}) = \left(\frac{n}{n-1} \widetilde{\sigma^2} \right) = \frac{n}{n-1} \left(\widetilde{\sigma^2} \right) = \frac{n}{n-1} \frac{n-1}{n} \sigma^2 = \sigma^2$$

We can call $\widetilde{\sigma^2}$ the **uncorrected sample variance**

Standard Error of Sample Mean

(Why divide by $n - 1$?)

- Only $n - 1$ independent pieces of information in $\{X_i - \bar{X}\}$ since $\sum_{i=1}^n (X_i - \bar{X}) = 0$
- Given $\{X_1 - \bar{X}, \dots, X_{i-1} - \bar{X}, X_{i+1} - \bar{X}, \dots, X_n - \bar{X}\}$, you can calculate $X_i - \bar{X}$
- you used one “degree-of-freedom” when you used the data to calculate \bar{X}
- If \bar{X} was obtained from a *different sample*, then you should divide by n , not $n - 1$, to get an unbiased estimator for σ^2

Standard Error of Sample Mean

For our data, we have

```
x <- dat$earn
N <- length(x)
muhat <- mean(x)
s2hat <- var(x)
muhatse <- sqrt(s2hat/N)
cat("sample mean:", round(muhat,3), " s.e.:", round(muhatse,3), "\n")
```

```
sample mean: 29.232  s.e.: 0.368
```

Consistency

$$E(\bar{X}) = \mu \text{ and } \text{Var}(\bar{X}) = \frac{\sigma^2}{n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

As $n \rightarrow \infty$, sample mean “converges” to μ

Convergence in Probability: A sequence of random variables $Y_n, n = 1, 2, \dots$, converges in probability to c if for any $\epsilon > 0$, we have

$$\lim_{n \rightarrow \infty} \Pr (|Y_n - c| \geq \epsilon) = 0 .$$

We say $Y_n \xrightarrow{p} c$

An estimator is **consistent** if it converges in probability to the true value of the parameter it is estimating

Consistency

Under our stated assumptions, the sample mean is consistent for the population mean

Khinchine's Weak Law of Large Numbers (WLLN) If $\{X_i\}_{i=1}^n$ is iid with $E(X_i) = \mu < \infty$ for all i , then

$$\bar{X}_n \xrightarrow{p} \mu$$

where \bar{X}_n is the sample mean based on n observations.

- There are many “Laws of Large Numbers” each stating different conditions under which the sample mean is consistent
- “Weak” refers to the kind of probabilistic convergence used here (there are others)
- Bias and variance going to zero is actually “convergence in mean square”, but this implies convergence in probability

Consistency (Simulation Example)

Suppose 200 people each took independent random samples of size n from population

Suppose population is well-represented by Chi-Sq(1) distribution (mean = 1)

Plot distribution of the 200 sample means for sample sizes

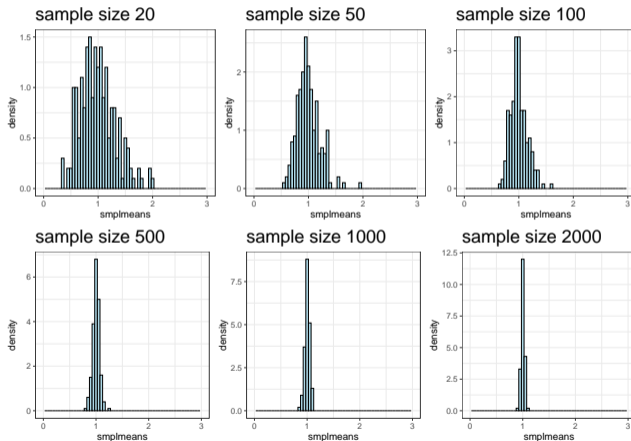
$n = 20, 50, 100, 500, 1000, 2000$

- distribution of samples means if 200 people drew sample of size 20
- distribution of samples means if 200 people drew sample of size 50

etc.

Plots demonstrate unbiasedness and consistency

Consistency (Simulation Example)



Consistency

An important property of convergence in probability: if $g(\cdot)$ is continuous, and $X_n \xrightarrow{p} c$, then $g(X_n) \rightarrow g(c)$

Suppose we want to estimate μ^2 . Use $\widehat{\mu^2} = \overline{X}_n^2$?

- \overline{X}_n^2 is **not** an unbiased estimator of μ^2 , since

$$\text{Var}(\overline{X}_n) = E(\overline{X}_n^2) - E(\overline{X}_n)^2 = E(\overline{X}_n^2) - \mu^2 \Rightarrow E(\overline{X}_n^2) = \mu^2 + \text{Var}(\overline{X}_n) > \mu^2$$

- \overline{X}_n^2 is a consistent estimator of μ^2 , since

$$\overline{X}_n \xrightarrow{p} \mu \Rightarrow \overline{X}_n^2 \xrightarrow{p} \mu^2$$

Consistency

$$\widetilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \text{ is consistent for } \sigma^2$$

Proof:

- X_i iid with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2 \Rightarrow X_i^2$ iid with $E(X_i^2) = \sigma^2 + \mu^2$
- $\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{p} \sigma^2 + \mu^2$ and $\bar{X}_n^2 \xrightarrow{p} \mu^2$
- Therefore $\widetilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \xrightarrow{p} \sigma^2 + \mu^2 - \mu^2 = \sigma^2$

$$\widehat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \text{ is also consistent for } \sigma^2 \text{ since } \widehat{\sigma}^2 = \underbrace{\frac{n}{n-1}}_{\rightarrow 1 \text{ as } n \rightarrow \infty} \widetilde{\sigma}^2$$

Hypothesis Testing (Two-Sided)

Suppose we want to test

$$H_0 : \mu = \mu_0 \text{ vs } H_A : \mu \neq \mu_0$$

Intuitive Idea:

- If $\mu = \mu_0$ we expect $\hat{\mu}$ to be “near” μ_0
- If $\hat{\mu}$ is far from μ_0 , perhaps $H_0 : \mu = \mu_0$ is incorrect
- If $\hat{\mu}$ is “too far” from μ_0 , take this as statistical evidence that $\mu \neq \mu_0$

But how far is too far?

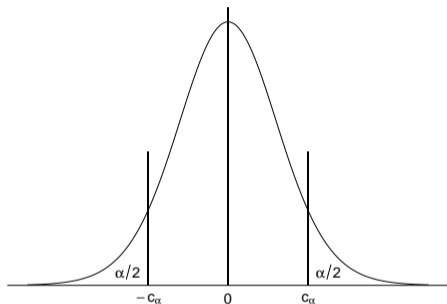
Hypothesis Testing (Two-Sided)

Assume for the moment that $X_i \stackrel{iid}{\sim} \text{Normal}(\mu_0, \sigma^2)$, $i = 1, \dots, n$

We have

$$\begin{aligned} X_i \stackrel{iid}{\sim} \text{Normal}(\mu_0, \sigma^2) &\implies \bar{X}_n \sim \text{Normal}\left(\mu_0, \frac{\sigma^2}{n}\right) \\ &\implies \frac{(\bar{X}_n - \mu_0)}{\sqrt{\sigma^2/n}} \sim \text{Normal}(0, 1) \\ &\implies \underbrace{\frac{(\bar{X}_n - \mu_0)}{\sqrt{\widehat{\sigma}^2/n}}}_{t\text{-statistic}} \sim t(n-1) \end{aligned}$$

Hypothesis Testing (Two-Sided)



Reject H_0 if $t > c_\alpha$ or $t < -c_\alpha$, where c_α is such that $\alpha = 0.01, 0.05, 0.10$ (“level of significance”)

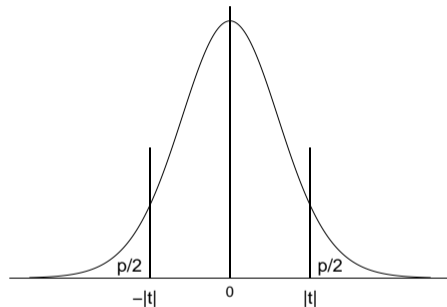
i.e., reject if $\Pr(|t| > c_\alpha) < \alpha$ given $\mu = \mu_0$ (Prob of rejecting correct null is α)

Hypothesis Testing (Two-Sided)

```
nVal <- c(20, 50, 100, 200, 400)
alphaVal <- c(0.01, 0.05, 0.1)
critVal <- matrix(rep(0,length(nVal)*length(alphaVal)), ncol = length(nVal))
colnames(critVal) <- paste0("n=",nVal)
rownames(critVal) <- paste0("alpha=",alphaVal)
for (i in 1:length(alphaVal)){
  for (j in 1:length(nVal)){
    critVal[i, j] = qt(1-alphaVal[i]/2, df=nVal[j]-1)
  }
}
round(critVal,3)
```

	n=20	n=50	n=100	n=200	n=400
alpha=0.01	2.861	2.680	2.626	2.601	2.588
alpha=0.05	2.093	2.010	1.984	1.972	1.966
alpha=0.1	1.729	1.677	1.660	1.653	1.649

Hypothesis Testing (Two-Sided)



Equivalently, reject $H_0 : \mu = \mu_0$ if “p-value” p in figure above is less than α , i.e., if the probability of obtaining a test statistic more extreme than the observed t -stat is less than α

Asymptotic Normality

When $n \rightarrow \infty$, the t-distribution converges to the Normal(0,1)

Then critical values $c_{0.01}$, $c_{0.05}$ and $c_{0.10}$ are 2.576, 1.96 and 1.645 respectively

- What if X_i does not have normal distribution? Then t -statistic does not have t distribution.

However, we have the following result

Lindeberg-Levy Central Limit Theorem: If $\{X_i\}_{i=1}^n$ are iid with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2 < \infty$ for all i , then

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \text{Normal}(0, \sigma^2)$$

Asymptotic Normality (Simulation Example)

Note: cannot talk about the distribution of \bar{X}_n as $n \rightarrow \infty$ since distribution collapses into single point

- Subtract μ to center it
- Multiply by \sqrt{n} so that variance does not go to zero

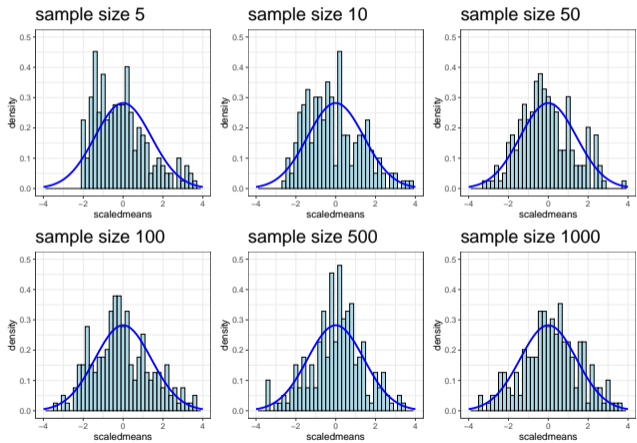
We can talk about distribution of $\sqrt{n}(\bar{X}_n - \mu)$ at $n \rightarrow \infty$

Continuation of simulation example of 200 people drawing independent samples from $\chi^2(1)$ population with sample sizes

$n = 5, 10, 50, 100, 500, 1000$

Plot distribution of $\sqrt{n}(\bar{X}_n - \mu)$ (here $\mu = 1$)

Asymptotic Normality (Simulation Example)



Hypothesis Testing (Two-Sided)

- “ \xrightarrow{d} ” means **convergence in distribution**
- when n is large, pdf of LHS is approximately the pdf of the Standard Normal
- Can also be shown that

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sqrt{\widehat{\sigma}^2}} = \frac{\bar{X}_n - \mu}{\sqrt{\widehat{\sigma}^2/n}} \xrightarrow{d} \text{Normal}(0, 1)$$

You can replace $\widehat{\sigma}^2$ with $\widetilde{\sigma}^2$ or any other consistent estimator of σ^2

That is, when n is large, can make the approximation $t \overset{a}{\sim} \text{Normal}(0, 1)$, where $\overset{a}{\sim}$ means “approximately distributed”, even when X_i is not normally distributed

Hypothesis Testing (Two-Sided) Example

Suppose we want to test

$$H_0 : \mu = 30 \text{ vs } H_A : \mu \neq 30$$

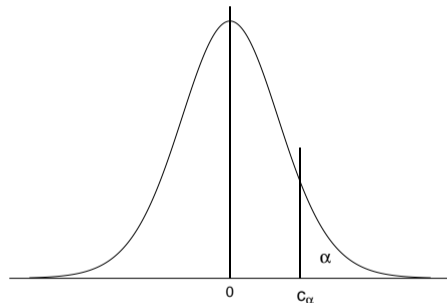
```
x <- dat$earn; n <- length(x); muhat <- mean(x); s2hat <- var(x)
tstat <- (muhat - 30)/sqrt(s2hat/n)
pval_t <- 2*pt(abs(tstat), df=n-1, lower.tail = FALSE)
pval_n <- 2*pnorm(abs(tstat), lower.tail = FALSE)
cat("t-stat:", tstat)
cat("\n p-value (t-dist):", pval_t)
cat("\n p-value (Standard Normal):", pval_n)
```

```
t-stat: -2.086885
 p-value (t-dist): 0.0369496
 p-value (Standard Normal): 0.03689851
```

Reject at 0.05 significance level but not 0.01

Hypothesis Testing (One-Sided)

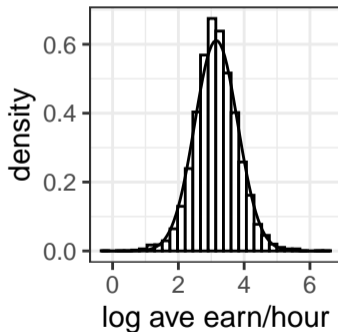
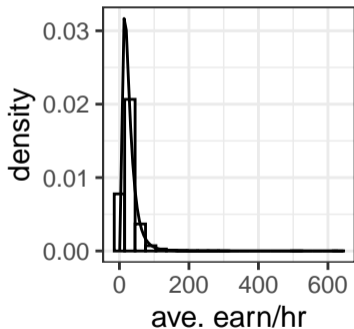
e.g., $H_0 : \mu \leq \mu_0$ vs $H_A : \mu > \mu_0$



Reject μ_0 if t-statistic is greater than c_α where c_α is that value such that $\Pr(t > c_\alpha) = \alpha$ under the null, $\alpha = 0.01, 0.05, 0.10$.

Estimation Again

Should we have worked with $\log(\text{earn})$ instead of earn ?



Estimation Again

Sample mean of *earn* is unbiased and consistent for population mean (no problem in this regard)

Finite sample distribution of sample mean may not be normal or even symmetric in smaller samples because of high skew in distribution of X_i

- Concerns about hypothesis tests
- Harder to interpret standard errors

Better to work with $\ln \text{earn}$? To estimate $E(\text{earn})$, probably not

When we get to regression, it will be easier and more sensible to work with $\ln \text{earn}$ instead of *earn*

Estimation Again

If you *do* work with $\ln \textit{earn}$ with intention to convert results to \textit{earn} , how do you do it?

To help us think about this, let's assume

$$\ln X_i \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2) \text{ for all } i$$

where X_i is earnings of individual i (seems reasonable!)

Then $X_i \stackrel{iid}{\sim} \text{Log-normal}(\mu, \sigma^2)$ for all i

- $E(X_i) = e^{\mu + \frac{1}{2}\sigma^2} = e^\mu e^{\frac{1}{2}\sigma^2}$
- $\text{Median}(X_i) = e^\mu$

Estimation Again

Can estimate $\mu = E(\ln X)$ and $\sigma^2 = \text{Var}(\ln X)$ in the usual way

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \ln X_i \quad \text{and} \quad \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (\ln X_i - \hat{\mu})^2$$

To convert mean and variance of $\ln X$ back to mean and variance of X

- estimate of mean hourly earnings: $e^{\hat{\mu}} e^{\frac{1}{2}\hat{\sigma}^2}$ (Not $e^{\hat{\mu}}$)
- estimate of median hourly earnings: $e^{\hat{\mu}}$

Also need to compute s.e. (use bootstrap?)

Estimation Again

Rough description of bootstrap idea

- Imagine again $R = 200$ people ($r = 1, \dots, 200$) each collecting individual size n samples from population, and calculated sample means for their own samples $\bar{X}_{n,r}$
- You can calculate the standard error for the sample mean using

$$\frac{1}{R-1} \sum_{r=1}^R (\bar{X}_{n,r} - \overline{\bar{X}_n})^2$$

where $\overline{\bar{X}_n}$ is the sample mean of the $R = 200$ sample means

- You can similarly get the standard error for the sample median

Estimation Again

Bootstrap idea: Create $R = 200$ "bootstrap" samples by re-sampling with replacement from your sample, i.e., from $\{X_1, \dots, X_n\}$, get

$$\{X_1^{(1)}, \dots, X_n^{(1)}\} \rightarrow \bar{X}_{n,1}, \text{Med}(X)_{n,1}$$

$$\vdots$$

$$\{X_1^{(r)}, \dots, X_n^{(r)}\} \rightarrow \bar{X}_{n,r}, \text{Med}(X)_{n,r}$$

$$\vdots$$

$$\{X_1^{(R)}, \dots, X_n^{(R)}\} \rightarrow \bar{X}_{n,R}, \text{Med}(X)_{n,R}$$

Calculate standard error of $\{\bar{X}_{n,r}\}_{r=1}^R$ and $\{\text{Med}(X)_{n,r}\}_{r=1}^R$ as standard error of your sample mean and sample median

Estimation Again

We apply this idea to get s.e. of mean and median of *earn* after mean and variance of $\ln \textit{earn}$

```
x <- log(dat$earn)
ln_mu_2_mu <- function(m, v){exp(m+0.5*v)}
ln_mu_2_md <- function(m, v){exp(m)}
m <- mean(x)
v <- var(x)
earnmean <- ln_mu_2_mu(m,v)
earnmed <- ln_mu_2_md(m,v)
set.seed(456)
R <- 200 ## Bootstrap replication sample
bvars <- bmeans <- bmeds <- rep(NA, R) ## To store the bootstrapped statistics
for (r in 1:R){
  xsmpb <- sample(x, 4946, replace=T) # Sample with replacement from orig. smp.
  m1 <- mean(xsmpb) # mean of bootstrap sample of ln(earn)
  v1 <- var(xsmpb) # variance of bootstrap sample of ln(earn)
  bmeans[r] <- ln_mu_2_mu(m1,v1) # convert to mean of earn, and store
  bmeds[r] <- ln_mu_2_md(m1,v1) # convert to median of earn, and store
}
```

