

# Describing Numbers

Paul E. Johnson<sup>1</sup> <sup>2</sup>

<sup>1</sup>Department of Political Science  
University of Kansas

<sup>2</sup>Center for Research Methods and Data Analysis, University of Kansas

2013

# What is this Presentation?

- A Brief summary of the idea “variable”
- Ways to Describe Numeric Variables
- Central Tendency: Mean, Median, Mode
- Dispersion: Variance, Standard Deviation, etc.
- Rescalings

# Variable

**Variable** a collection of scores that represent observations.

Example:

$$height = \{6.0, 5.1, 4.2, 5.8, 5.4\} \quad (1)$$

**Subscript  $height_i$ :**  $height_1$  is observation 1,  $height_2$  is observation 2, and so forth

# Common Notation

More abstractly

$$x = \{x_1, x_2, x_3, x_4, \dots, x_N\} \quad (2)$$

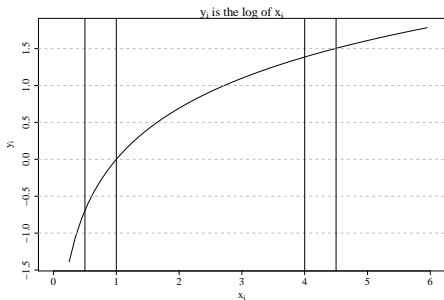
Or perhaps more succinctly

$$x_i, \text{ for } i \in \{1, \dots, N\} \quad (3)$$

- 1**  $N$ : capital  $N$  refers to “sample size” or “number of observations” (in most social sciences).
- 2** Usually, when I talk about  $x_i$ , I mean to refer to any of the individual observations in  $x$ .
- 3** Set notation
  - 1**  $\in$  means “element of,” as in  $i \in N$  or  $x_2 \in X = \{x_1, x_2, \dots, x_N\}$ .
  - 2**  $\forall$  abbreviation of “for all”, so “for  $i \in N$ ” might be  $\forall i \in N$ .

# Numeric Variables

- NUMERIC variables: accept mathematical transformations
- The range from  $\{minimum, maximum\}$  is (subjectively) meaningful
- From  $x_i$  to  $2 \times x_i$ : there is twice as much of it
- Analysis may be altered (improved or damaged) by transformations



- $\log()$  magnifies the importance of a step from 0.5 to 1 and shrinks the importance of a step from 4 to 4.5.

# Here's your Mission

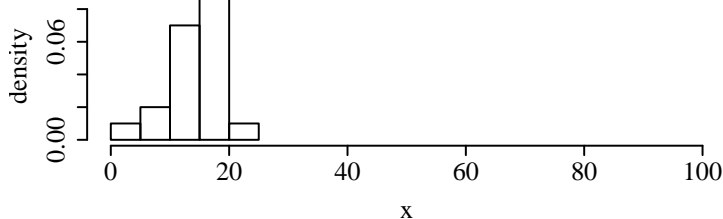
- You found some scores (data)
- Can your data for yourself
- Tell/show other people about it
  - You need terminology: describe
  - Show a picture
  - Summarize with numbers.

## Terminology to Describe Variables

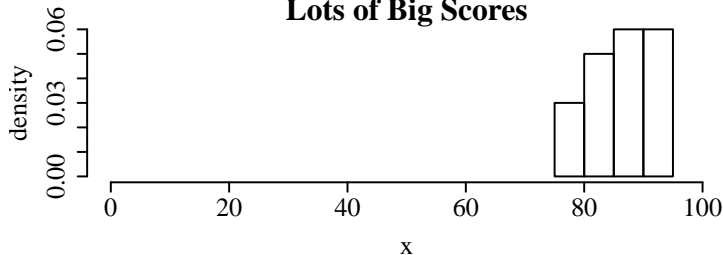
- Central Tendency: Where, “generally” are the scores? Is there a “meaningful” (subjective) characterization of where “most” scores are situated
- Dispersion: How “spread out” are the scores? Is it not meaningful to talk about a “typical” observation?
- Shape of Distributions: Do the observations appear to be
  - Unimodal (one most-likely score, others less likely)
  - Symmetric or Skewed

# Histograms: Compare Two Variables

## Lots of Low Scores



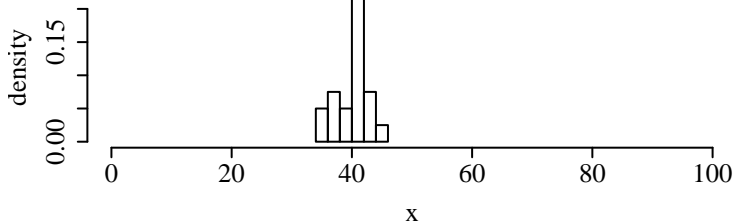
## Lots of Big Scores



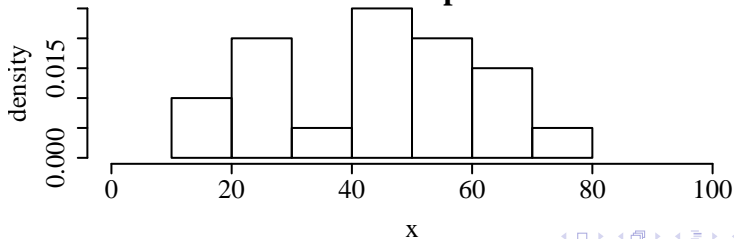


# Histograms: Compare Two Variables

**These are all clumped up**



**These are all spread out**

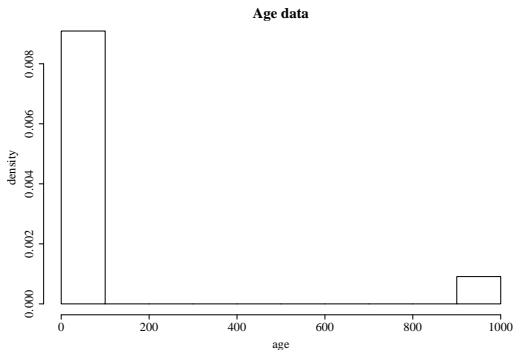


## Define "Histogram", Please

- Group observations into "bins" of similar scores
- Draw bars to represent the proportion of all scores that fall into each bin
- The areas of the bars should sum to 1.0

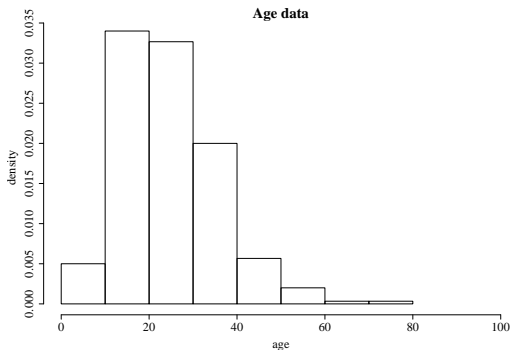
# Histograms: Check for Data Errors

- Suppose your data is supposed to be human age
- But somebody put in 999 for “missing” data points



# Histograms: Check for Data Errors

- If you ignore (remove) the cases that are equal to 999 (or set them to NA)
- Generally, whenever you get new data from anybody/anywhere, a histogram is a good “first check” on it.



## Various "transformations" might be applied

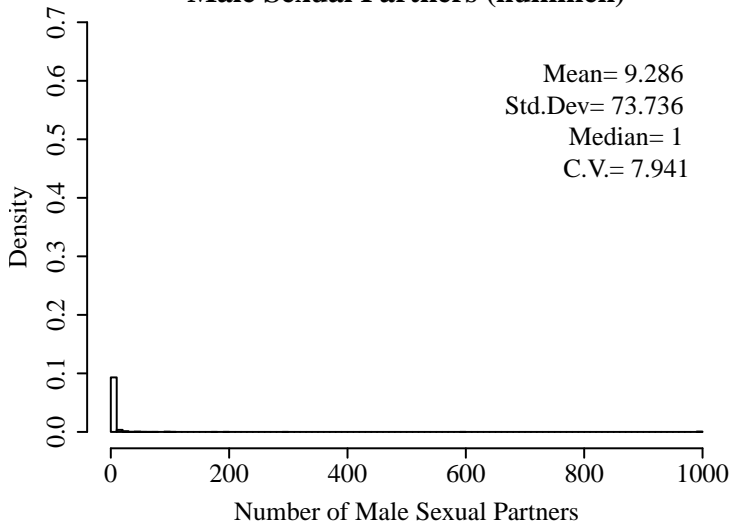
- Fox, Ch. 2 reviews various ways to re-scale data to make it 'fit' some statistical tests
- I'm cautious about fiddling with data
- Some transformations are not "harmless"
- Goal: Be honest with self & others about changes applied to data, including
  - omission of missing or extreme observations
  - multiplicative re-scaling
  - nonlinear transformations (log, Box-Cox, etc.)

# Some Examples from the General Social Survey

`/stat/DataSets/GSS/gss-subset2.Rda`

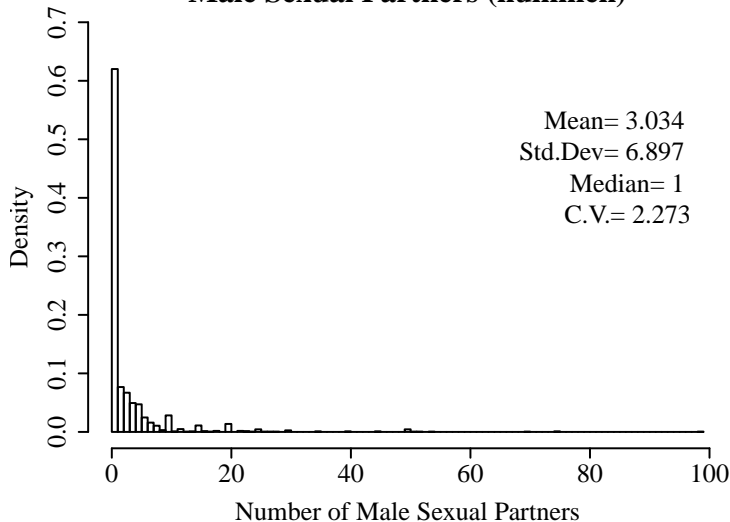
# Histogram: Spot Typos/Unusual Scores

## Male Sexual Partners (nummen)



# Histogram: Eliminate values greater than 99

## Male Sexual Partners (nummen)





## The Size of the Bins Can Make a Difference

- Narrow bars have more detail, possibly less generalizability (harder to see patterns)
- Wide bars smooth out too many bumps, hide details
- Many algorithms proposed to choose bin width to automate production of “good” histograms.

# Histogram: Fatter Bars!

## Male Sexual Partners (nummen)

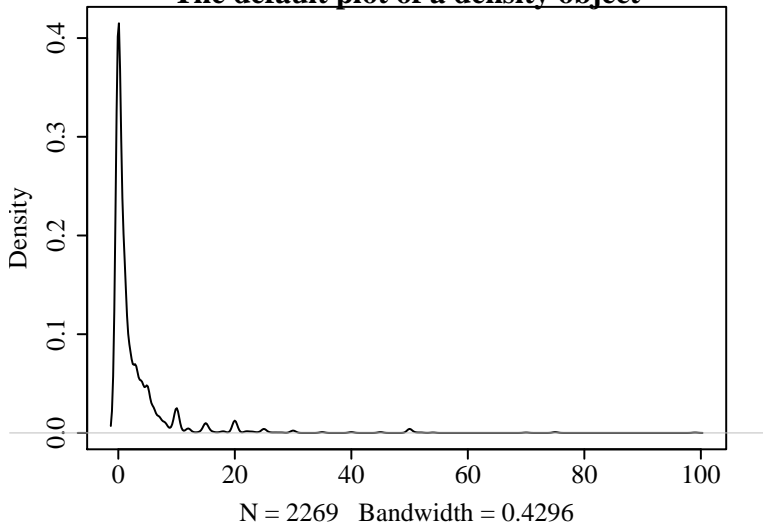


## A Smoothing Curve: Kernel Density Estimate (KDE)

- Because of the (subjective) “bin width” problem, other density estimation methods have been developed
- The kernel density estimate is a “smoothing” method that estimates the density at each value, putting more weight on nearby observations than far away ones.
- Some propose to replace histograms with KDE

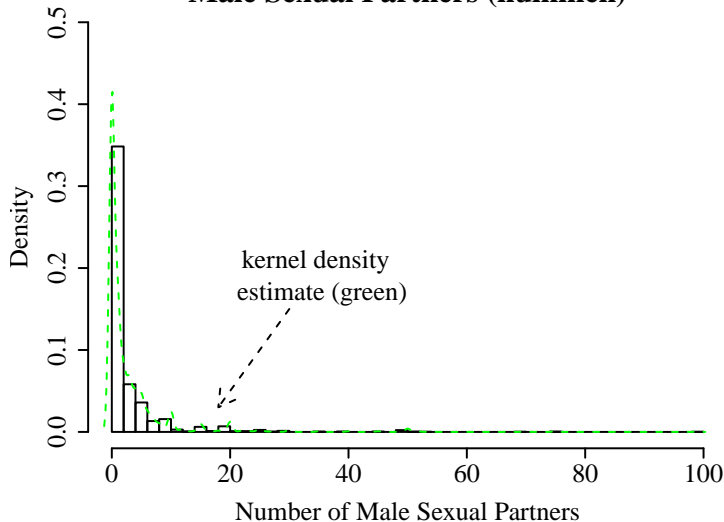
# The Density Estimates

**The default plot of a density object**



# Histogram with Density Super-imposed

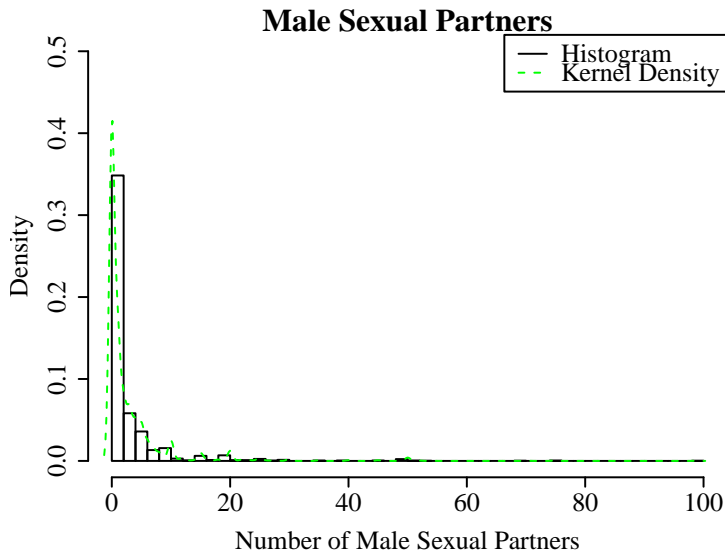
## Male Sexual Partners (nummen)



## Histogram: More on Customizing Histograms

- My lectures in guides/Rcourse (plot-1, plot-2) have plenty of additional detail on beautifying plots.

# Histogram: with a "legend"



## Convey Same Info Without Graph?

- What if your publisher will not allow you the space for a histogram?
- Convey same information without a picture?
- Need to develop terminology to describe and compare what we see.



# Mean = Average

- The mean is a widely-used indicator of a distribution's central tendency

**Mean:** (same meaning as “average”).

$$\bar{x} = \frac{\text{Sum Of All Scores}}{\text{Number Of Scores}} \quad (4)$$

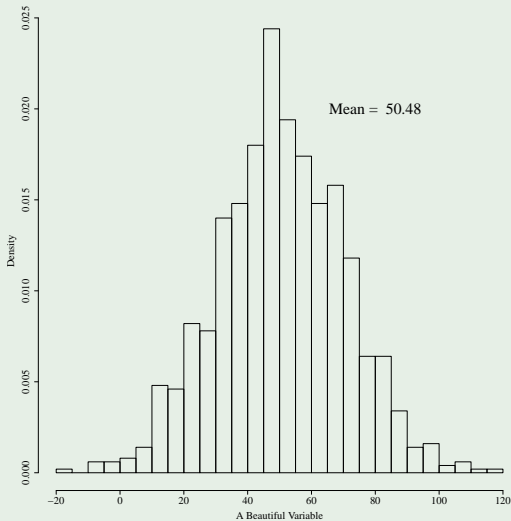
Suppose  $x$  is a variable. Add up the scores, then divide by  $N$ .

$$\bar{x} = \frac{\sum_{i=1}^N x_i}{N} \quad (5)$$

- 1 About Notation:  $\bar{x}$  is common notation, but not required. Sometimes I write  $Mean(x)$  or  $\hat{\mu}$ , depending on the context.

- I manufactured a sample of pleasantly symmetrical random data
- The sample mean is 50.485
- Appears (to me)
  - unimodal (one peak)
  - symmetric (more or less)

## A Histogram with 30 Bins



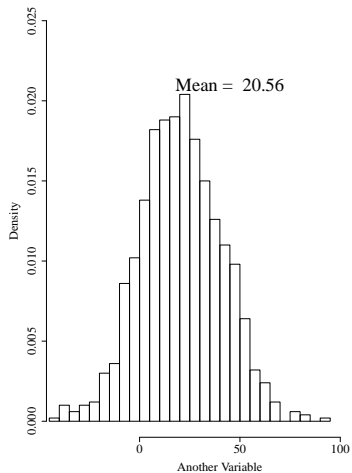
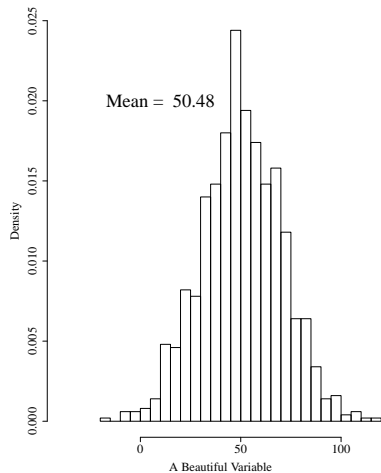
## You too can manufacture Normal samples

- I used R's `rnorm` function to draw some example observations

```
set.seed(1234321)
myx <- rnorm(1000, mean=50, sd=20)
```

- That creates 1000 observations from the Normal distribution,  $N(50, 20^2)$
- We specify 2 parameters
  - 50 is the parameter mu ( $\mu$ ), the "true mean"
  - 20 is the parameter sigma ( $\sigma$ ), which controls the "dispersion" of the scores.
- "Gaussian distribution" another name for the Normal.
- In case you wondered, the sample standard deviation is 19.977

# Compare 2 variables



# Variance

**Variance:** the average of squared deviations about the mean.

Calculate the difference between the  $i$ 'th case and the mean:

$$x_i - \bar{x} \quad (6)$$

Square that:

$$(x_i - \bar{x})^2 \quad (7)$$

Do the same for all and add them up:

$$\sum_{i=1}^N (x_i - \bar{x})^2 \quad (8)$$

Then divide by  $N$ .

$$\text{Var}(x) = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N} \quad (9)$$

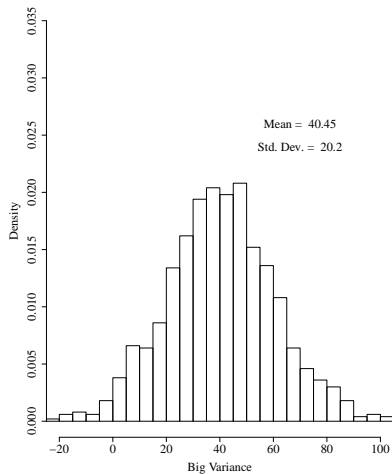
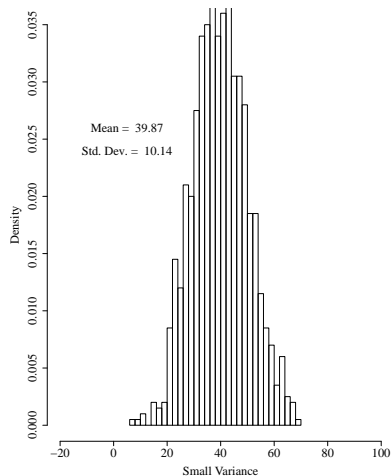
# Standard Deviation

**Standard Deviation:** the square root of the variance.

$$\text{Std.Dev.}(x) = \sqrt{\text{Var}(x)} = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}} \quad (10)$$

- Var and Std.Dev. serve same purpose.
- Std.Dev. has an advantage: it is measured (roughly speaking) on the same scale as the mean. (see below on “scaling”)
- Please don't worry right now about the need to divide by  $N - 1$  instead of  $N$ . That distraction is not needed at this stage.

# Compare 2 variables

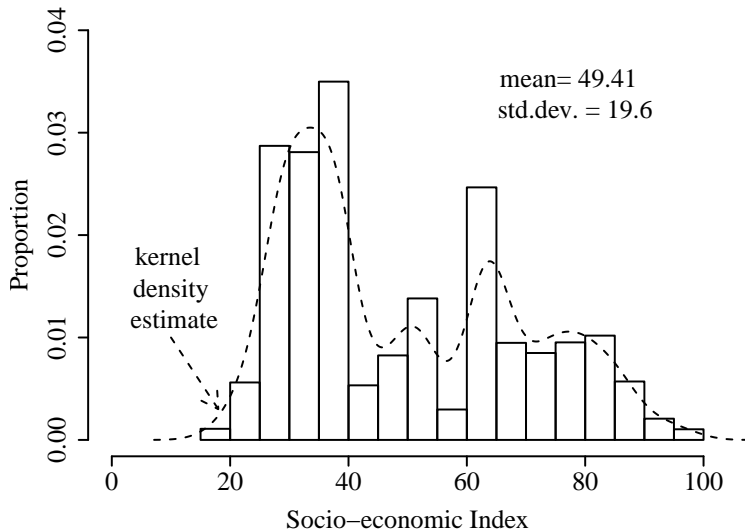


## About Notation

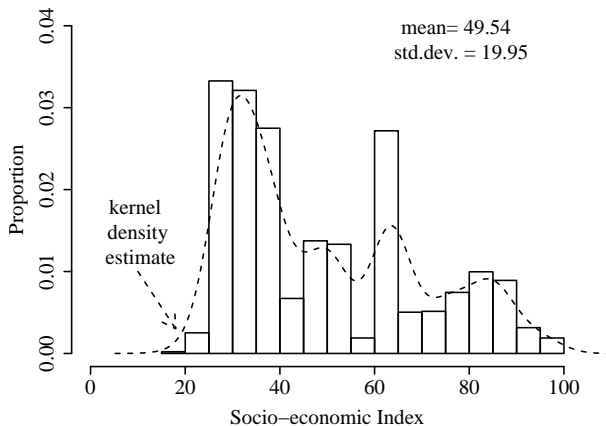
- 1 Some traditional stats books call the observed variance  $s^2$  and the “true variance” (of which it is an estimate)  $\sigma^2$ . Some books say “True standard deviation” is  $\sigma$  and estimate is  $s$ .
- 2 Some books call estimated variance  $\widehat{\sigma^2}$ . I like that because I don't need separate letters  $\sigma$  and  $s$ .
- 3 Some books call the “true” variance  $Var(x)$  while an estimate is  $\widehat{Var(x)}$ . I like that too, its easy to remember what's what.



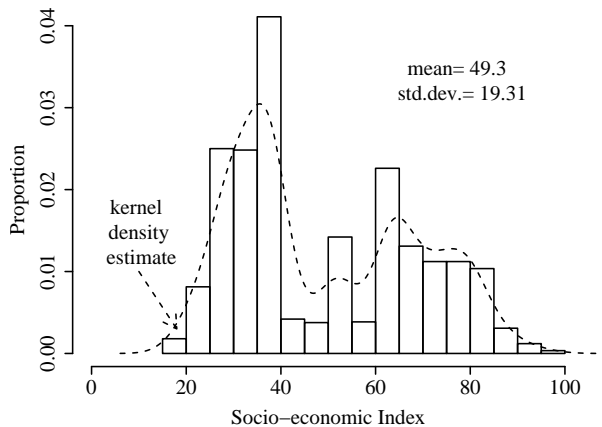
# Socio Economic Status



# Socio Economic Status: Only Men



# Socio Economic Status: Women



## Other Diversity Indicators

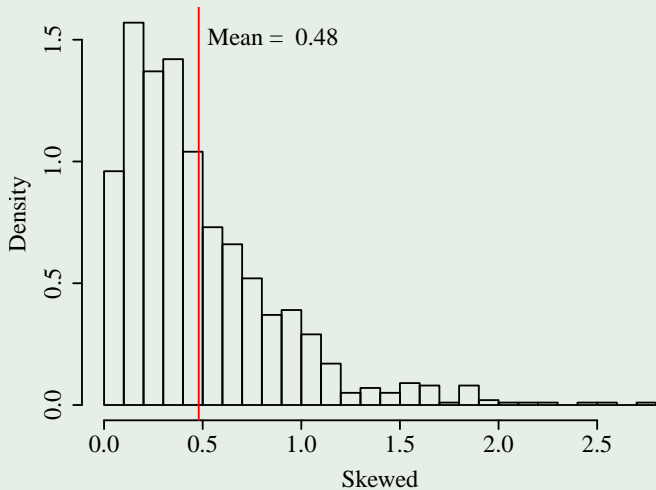
- Inter-Quartile range: group data by ordered quarters, and then think of the range between 25 percentile and 75 percentile as a diversity indicator.
- Many possible diversity indicators, including
  - gini index (often used for income inequality)
  - the mean of absolute valued differences

$$\text{Mean Absolute Deviation} = \frac{\sum_{i=1}^N |x - \bar{x}|}{N} \quad (11)$$

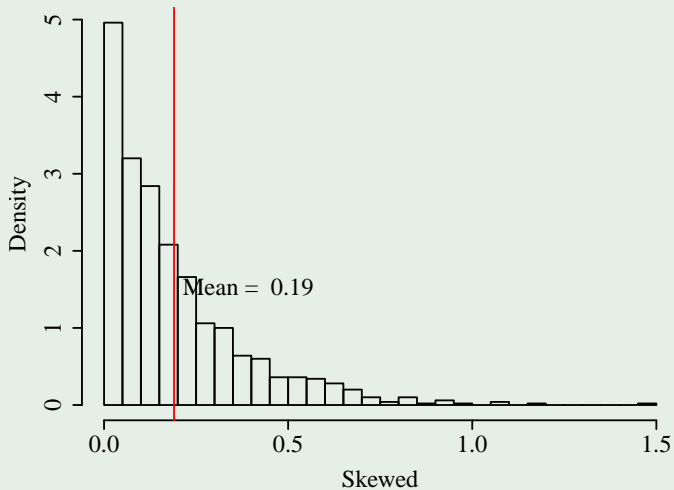
# Symmetry Definition

- A distribution is symmetric if the chance of observing a score  $\bar{x} - c$  is the same as observing  $\bar{x} + c$ .
- If a distribution is symmetric, then we have no trouble conveying the idea of its 'location'.
- The mean is in the middle!

## A Nonsymmetric Distribution



## Another Nonsymmetric Distribution



# Median: Center Case

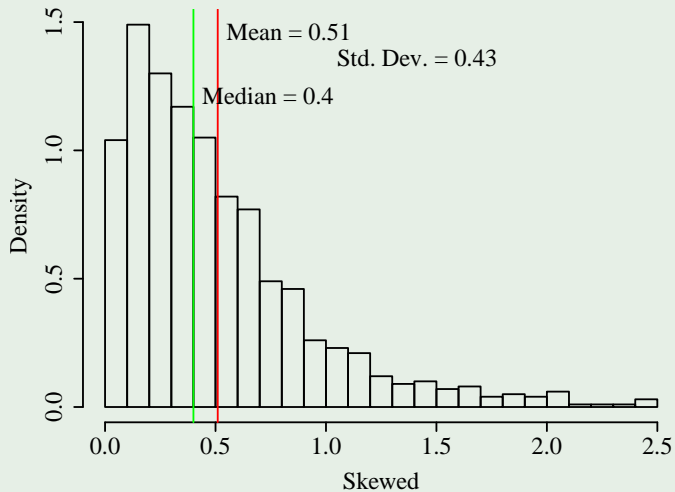
**Median:** The “center observation,” the number of observations that are larger equals the number that is smaller.

Questions:

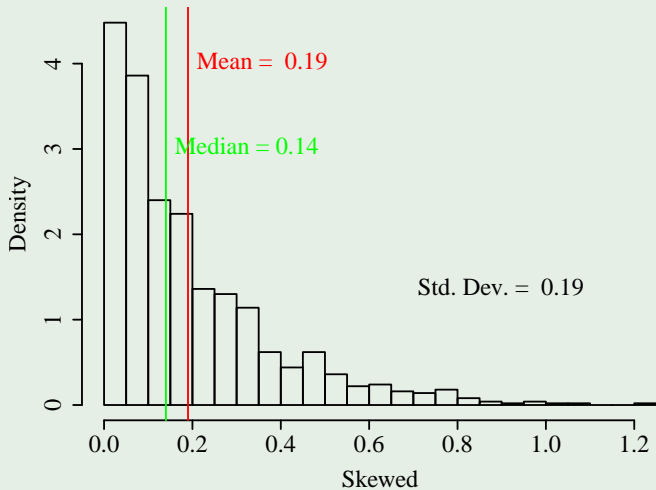
- 1 When do you think the mean and median are likely to be the same?
- 2 Can you think of a situation in which the median may be more meaningful than the mean?



## Add the median. Helpful?



## Another Nonsymmetric Distribution





## Note about Level of Measurement

- Mean only useful if we have numerical data (silly to average “low”, “medium”, “high”)
- Median requires ordered data, either numerical or ordered categorical
- Problem with the mean: it is distorted by a change in one value on either side (change one 50 to 5,000,000 and note the mean changes)
- Median is a more “robust” estimate (jargon: high ‘breakdown point’)

## Should the Scale Matter?

- The temperature in Celsius is 10. The temperature in Fahrenheit is 50 ( $32+9/5*10$ ).
- My income in dollars is 68,000. My income in Euros is 43,000 and in Pesos it is 1,126,123.
- Sometimes, we receive data in one format, but convert to another
- Simple scale conversions SHOULD NOT substantively change the conclusions we will draw.
- If simple scale conversions seem to matter, be VERY cautious.

## The Mean Scales With The Data

- Take variable  $X = \{x_1, x_2, \dots, x_N\}$ , and multiply each value by 10 to create  $newX$

$$newX = \{10x_1, 10x_2, \dots, 10x_N\} \quad (12)$$

- The mean of  $newX$  is obviously 10 the mean of old  $x$ . See?

$$Mean(newX) = \frac{10x_1 + 10x_2 + \dots + 10x_N}{N} = 10 \frac{\sum_{i=1}^N x_i}{N}$$

$$Mean(newX) = \overline{newX} = 10 \times \bar{x}$$

- Generally (meaning always), the mean of  $(k \times X)$  is equal to  $k$  times the mean of  $X$ .

## My First Big Fact

- State that as a theorem.  $k_1$  and  $k_2$  are any non-zero constants.  $X$  is any variable. Create a new variable  $newX = k_1 + k_2X$

The Mean scales proportionally. Given constants  $k_1, k_2$

$$Mean(k_1 + k_2X) = k_1 + k_2 \times Mean(X) \quad (13)$$

- The point: The Mean changes in a completely predictable way when the data is re-scaled by addition and multiplication. Just apply same same re-scaling to the old mean.

## The Variance Doesn't Scale Proportionally

- Suppose variance of  $X$  is  $\text{var}(X)$
- Create  $\text{new}X$  by multiplying by 10,  $\text{new}X = 10 \cdot X$
- The variance of  $\text{new}X$  is  $10^2 \text{Var}(X)$



## General Result for Variance of Re-scaled Variables

Calculate the Variance of a re-scaled Variable,  $X$ . Given  $k_1$  ,  $k_2$

$$\text{Var}(k_1 + k_2 \cdot X) = k_2^2 \cdot \text{Var}(X) \quad (14)$$

- Adding  $k_1$  does not change the dispersion at all, it just shifts the scores.
- The variance of  $\text{new}X = k_1 + k_2X$  is  $k_2^2 \times \text{Var}(x)$

Implication: Don't re-calculate mean and variance if  $x$  is proportionally re-scaled.

- Celsius temperature data,  $x$ . Suppose the mean is, 100.
- Rescale that data to Fahrenheit

$$xF_i = 32 + \frac{9}{5}x_i \quad (15)$$

- Some students want to re-run  $xF_i$  through the mean function, but they don't need to.
- The mean of  $xF$  is  
 $32 + (9/5)Mean(x) = 32 + (9/5)100 = 212.$

## But the Standard Deviation Scales Proportionally!

- The variance of  $xF$  is  $(9/5)^2 \times \text{Var}(x)$ , which is NOT linear
- However, recall standard deviation is  $\sqrt{\text{Var}(x)}$ , so the standard deviation would be

$$\text{Std.Dev.}(xF) = \sqrt{(9/5)^2 \times \text{Var}(x)} = (9/5) \times \text{Std.Dev.}(x) \quad (16)$$

- Like the mean, the standard deviation scales proportionally.

Standard Deviation of  $kX$  is  $k \times \text{Std.Dev.}(X)$

$$\text{Std.Dev}(k \cdot X) = k \cdot \text{Std.Dev}(X) \quad (17)$$

## The ratio $Std.Dev./Mean$ is Also Scale Invariant

- Recall  $Mean(k \cdot x) = kMean(x)$
- And  $Std.Dev.(k \cdot x) = kStd.Dev.(x)$
- Then the ratio of the mean to the standard deviation is not affected by  $k$

$$\frac{Mean(k \cdot x)}{Std.Dev.(k \cdot x)} = \frac{k \cdot Mean(x)}{k \cdot Std.Dev.(x)} = \frac{Mean(x)}{Std.Dev.(x)} \quad (18)$$

- And the converse is also true

$$\frac{k \cdot Std.Dev(x)}{k \cdot Mean(x)} = \frac{Std.Dev(x)}{Mean(x)} \quad (19)$$

# Coefficient of Variation is $\text{Std.Dev}(x)/M(x)$

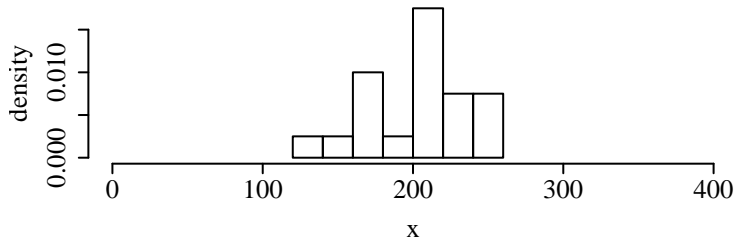
Coefficient of variation, CV.

Question: is “this distribution” more “spread out” than “that one”?

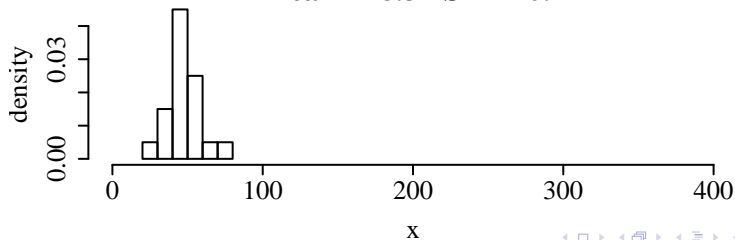
- This is a difficult, possibly silly question when distributions are fundamentally different
- But, if they have roughly the same “shape”, then the re-scaling might make them comparable.

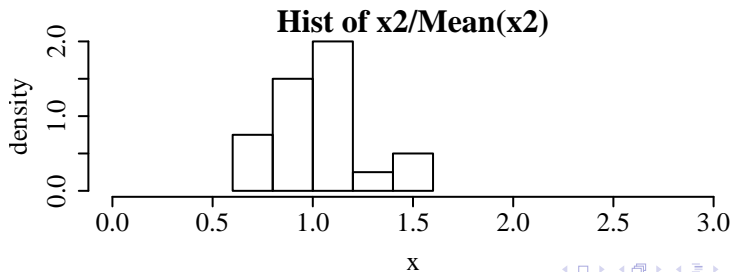
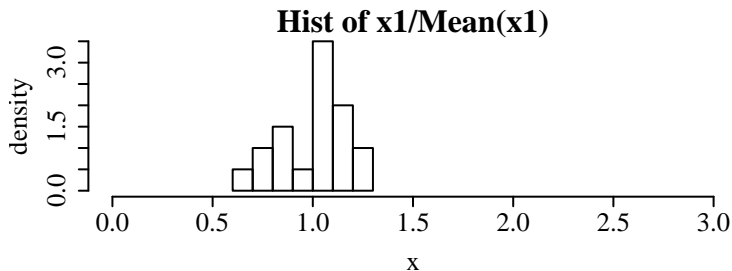
# Compare dispersion of 2 disparate variables

**Mean= 202.84 SD= 33.08**



**Mean= 48.32 SD= 10.4**



Compare 2: plot  $x/\text{Mean}(x)$ 

## Summarize those 2 variables

	top	bottom
Min.	129.9000000	29.6700000
1st Qu.	174.5000000	43.6800000
Median	214.0000000	48.9300000
Mean	202.8000000	48.3200000
3rd Qu.	225.8000000	50.6000000
Max.	246.5000000	73.1700000
mean	202.8352460	48.3206232
sd	33.0755854	10.3983793
sd.over.mean	0.1630663	0.2151955



## Mean-center $x_i$

- Mean centered data (aka “data in deviations form”)

$$\text{Mean Centered}(x_i) = x_i - \text{Mean}(x_i) \quad (20)$$

- Do we need abbreviation for that?  $x_i^{MC}$  or  $\tilde{x}_i$  or ?
- The mean of a centered variable is always 0
- The variance and standard deviation are unchanged by centering
- Sometimes mean-centered data is used to facilitate interpretation of results in some models.

## Standardized Data: special name for *Mean Centered*( $x_i$ )/*Std.Dev*( $x$ )

### Standardized Variables.

- Standardize means “divide *Mean Centered*( $x_i$ ) by standard deviation”.

$$\frac{x_i - \bar{x}}{\sigma_x} \quad (21)$$

- Since *M*( $x$ )/*Std.Dev*( $x$ ) is scale invariant, it makes *Mean Centered*( $x_i$ )/*Std.Dev*( $x$ ) will also be unaffected by re-scaling of the observations.
- The letter “Z” is often used to refer to standardized variables.

# Standardized implies Mean 0, Std.Dev 1



$$\text{Mean of } Z = \bar{Z} = 0$$

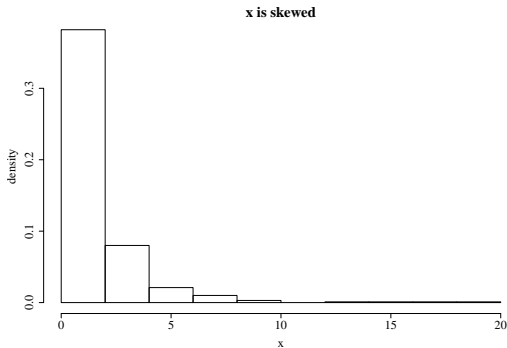


$$\text{Std.Dev.}(Z) = SD(Z) = 1$$

- It is strictly a matter of convention to standardize variables.
- Standardization may allow comparison of distributions, but it may not (more advanced problem I'm not willing to go into)

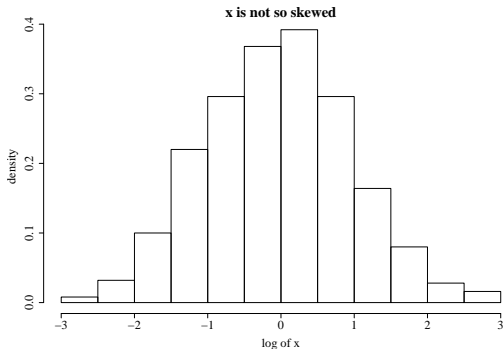
# The log is the most commonly applied nonlinear transformation

- We often gather data that is “clumped” on the left
- Examples, income, education



# The log is the most commonly applied nonlinear transformation

- The distribution of  $\log(x)$  appears more symmetric



## Difficult to say for sure if logging a variable is good or bad

- Some methods books will recommend logging all variables, claiming that it almost always makes analysis “work better” in some sense.
- Please just remember it is a possibility

## R functions to remember

`x` is a variable

- `mean(x, na.rm = TRUE)`
- `sd(x, na.rm = TRUE)`
- `var(x, na.rm = TRUE)`
- `median(x, na.rm = TRUE)`
- `range(x, na.rm = TRUE)`
- `quantile(x, na.rm = TRUE)`
- `summary(x)`
- `rockchalk::summarize(x)`
- `hist(x, prob = TRUE)`
- `xdens <- density(x)`
- `lines(xdens)`
- `plot(xdens)`

# Problems

- 1 Better run hist a few times. If you have R handy, try this

```
x1 <- rnorm(100, m=20, s=10)
hist(x1, prob = TRUE, main = "mean of 20, standard
      deviation of 10")
den1 <- density(x1)
lines(den1, col = "red", lty = 4)
plot(den1)
x2 <- x1
x2[98:100] <- 999
hist(x2, prob = TRUE, main = "mean of 20, standard
      deviation of 10")
den2 <- density(x2)
lines(den2, col = "red", lty = 4)
plot(den2)
```

- 2 There are some weird arguments you can use with hist. Try this.

- 1 Change prob = TRUE to prob = FALSE.



## Problems ...

- Change the number of bins by setting `breaks = 40` or `breaks = 4`.
- Fiddle with the appearance by adjusting `ylim = c(0, 1)` or `ylim = c(0.2, 0.8)`.
- Now make some weird looking data and do the same. This will create 300 funny looking observations in variable `x`, I'm pretty sure:

```
x1 <- rnorm(100, m=30, s=10)
x2 <- rpois(100, lambda=1)
x3 <- rnorm(100, m=80, s=20)
x <- c(x1, x2, x3)
```

- use the functions “mean”, “var” and “sd” on each one of those.
- Make a histogram for each.
- It may be beyond your R skills now to insert labels for the mean or standard deviation, but you could pencil them in if you had paper.

# Problems ...

- 4** Try this. It draws a fresh set of data

```
x1 <- rnorm(100, m = 30, s = 10)
x2 <- rpois(100, lambda = 1)
x3 <- rnorm(100, m = 80, s = 20)
dat <- data.frame(x1, x2, x3)
rm(x1, x2, x3)
library(rockchalk)
summarize(dat)
sapply(dat, mean)
var(dat)
```

- 5** Missing data is a sad fact of life.

- 1** Insert some missings (by setting cases to NA). Note how I keep it interesting by changing idioms for access to rows and columns.

```
dat$x1[c(1, 2, 55)] <- NA
dat[ c(56, 99), 2] <- NA
```

# Problems ...

- 2 Try to collect statistics about that data.
- 3 Understand why we need to rewrite the use of the mean function to tolerate NA?

```
mean(dat$x1)
mean(dat$x1, na.rm = TRUE)
sapply(dat, mean, na.rm = TRUE)
```