# Chapter 2

# Describing Variables

2.5   Measures of Dispersion

# Measures of Dispersion

Measures of dispersion indicate the amount of variation or "average differences" among the scores in a frequency distribution.

We're less familiar with such concepts in daily life, although a range of values is sometimes reported:

- **Today's forecast high temp will be 59-62 degrees**

- **N. Korea's Taepodong missile has a reported range of 2,400 to 3,600 miles**

- **Gallup Poll reported 51% of a national sample agree that President Obama is doing a good job, with a "margin of error" of $\pm 3\%$**

# Discrete Variable Dispersion Measures

**Index of Diversity (D)** measures whether two randomly selected observations are likely to fall into the same or different categories

$$D = 1 - \sum_{i=1}^{K} p_i^2$$

Higher D indicates the cases are <u>more equally spread</u> across a variable's *K* categories (i.e., they are less concentrated)

# Calculate D for these four GSS regions of residence:

| Region | $p_i$ | $(p_i)^2$ |
|--------|-------|-----------|
| NORTH EAST | .175 | _____ |
| MIDWEST | .215 | _____ |
| SOUTH | .361 | _____ |
| WEST | .248 | _____ |

$$\sum_{i=1}^{K} p_i^2 = \underline{\hspace{3cm}}$$

$$D = 1 - \Sigma p_i^2 = \underline{\hspace{5cm}}$$

The **Index of Qualitative Variation (IQV)** adjusts D for the number of categories, *K*

$$IQV = \frac{K}{K-1}(D)$$

IQV gives a bigger "boost" to *D* for a variable with fewer categories, thus allowing comparison of its dispersion to a variable that has more categories

Sally and three friends buy a 12-pack of beer (144 oz.). Ted and seven friends buy two 12-packs (288 oz.). Which distribution of beer is "fairer" (more equally distributed within each set of drinkers)?

**Sally: 20, 28, 46, 50 oz.**

**Ted:   20, 28, 32, 36, 40, 40, 44, 48 oz.**

$$IQV = \left(\frac{K}{K-1}\right)\left(1 - \sum_{i=1}^{K} p_i^2\right)$$
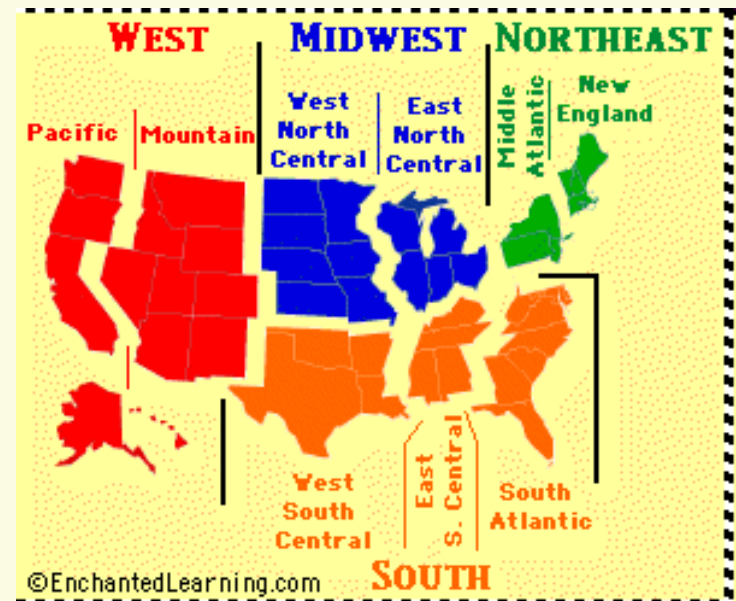
**IQV$_{Sally}$ =** _____

**IQV$_{Ted}$ =** _____

Indices of Diversity for proportions of U.S. population living in 4 Census regions and the distribution in 9 Census regions:

Four-region D = 0.731

Nine-region D = 0.855

The population seems more equally spread among the 9 regions than among the 4 regions. However, …



calculate the IQVs for both measures. Do these two population distributions now seem differently dispersed?

Four-region IQV = _____

Nine-region IQV = _____

**Range** the difference between largest and smallest scores in a <u>continuous variable</u> distribution

What are the ranges for these GSS variables?

|  | <u>Min.-Max.</u> | <u>Range</u> |
|---|---|---|
| EDUC: | 0 to 20 years | _____ |
| AGE: | 18 to 89 years | _____ |
| PRESTG80: | 17 to 86 points | _____ |
| PAPRES80: | 17 to 86 points | _____ |

# Average Absolute Deviation (AAD)

Read this subsection (pp. 48-49) for yourself, as background info for the variance & standard deviation

Because ADD is never used in research statistics, we won't spend any time on it in lecture

# Variance and Standard Deviation

Together with the mean, the variance (and its kin, the standard deviation) are the workhorse statistics for describing continuous variables

**Variance** the mean (average) squared deviation of a continuous distribution

The <u>deviation</u> ($d_i$) of case <u>i</u> is the difference between its score $Y_i$ and the distribution's mean:

$$d_i = Y_i - \overline{Y}$$

To calculate the variance of a <u>sample</u> of N cases:

- Compute and square each deviation

- Add them up

- Divide the sum by N - 1

$$s_Y^2 = \frac{\sum\limits_{i=1}^{N} (Y_i - \overline{Y})^2}{N-1} = \frac{\sum d_i^2}{N-1}$$

Reason for using N-1, not N, will be explained later.

**Standard deviation** the positive square root of the variance

This transformation avoids the unclear meaning of squared measurement units; e.g., years-squared

The standard deviation of a sample:

$$S_Y = \sqrt{S_Y^2}$$

# Calculate $s^2$ and $s$ for these 10 scores

$$Y_i \;-\; \overline{Y} \;=\; d_i \qquad (d_i)^2$$

| | | | |
|---|---|---|---|
| 2 - 2 = | _____ | _____ |
| 0 - 2 = | _____ | _____ |
| 4 - 2 = | _____ | _____ |
| 1 - 2 = | _____ | _____ |
| 6 - 2 = | _____ | _____ |
| 3 - 2 = | _____ | _____ |
| 1 - 2 = | _____ | _____ |
| 2 - 2 = | _____ | _____ |
| 1 - 2 = | _____ | _____ |
| 0 - 2 = | _____ | _____ |

$$\sum_{i=1}^{10} (d_i)^2 = \underline{\hspace{3cm}}$$

$$s_Y^2 = \sum (d_i)^2 \; / \; (N-1) =$$

$$\underline{\hspace{5cm}}$$

$$s_Y = \sqrt{s_Y^2} = \underline{\hspace{3cm}}$$

To calculate the variance of a <u>dichotomy</u>, just multiply both proportions: $$s_Y^2 = (p_0)(p_1)$$

The 2008 GSS asked, "Do you favor or oppose the death penalty for persons convicted of murder?" What is its variance?

| CAPPUN | $p_i$ |
| --- | --- |
| 1 FAVOR | .66 |
| 0 OPPOSE | .34 |

$$s_Y^2 = \rule{4cm}{0.4pt}$$

A item about having ever used crack cocaine was split more unevenly. Is its variance larger or smaller than CAPPUN's?

| EVCRACK | $p_i$ |
| --- | --- |
| 1 YES | .06 |
| 0 NO | .94 |

$$s_Y^2 = \rule{4cm}{0.4pt}$$

# Variance of a Grouped Frequency Distribution

Use the variance formula but multiply each squared deviation by its relative frequency ($f_i$), then sum the products across all *K* categories:

$$s_Y^2 = \frac{\sum_{i=1}^{K} (Y_i - \overline{Y})^2 (f_i)}{N-1} = \frac{\sum (d_i^2)(f_i)}{N-1}$$

## What is the variance of these grouped data?

HOMOSEX1  "What about sexual relations between two adults of the same sex; is it …"

[Mean = 2.15 for N = 1,309]

| Response | $Y_i$ | $f_i$ | $(d_i)^2(f_i)$ |
|---|---|---|---|
| Always wrong | 1 | 733 | _____ |
| Almost always | 2 | 67 | _____ |
| Only sometimes | 3 | 88 | _____ |
| Not wrong at all | 4 | 421 | _____ |

$$s_Y^2 = \frac{\sum_{i=1}^{K}(d_i)^2(f_i)}{N-1} = \underline{\hspace{4cm}}$$

**Skewness** describes nonsymmetry (lack of a mirror-image) in a continuous distribution

It compares the mean and the median:

$$\text{Skewness} = \frac{3(\overline{Y} - \text{Mdn})}{S_Y}$$

- Positive skew has a "tail" to right of Mdn
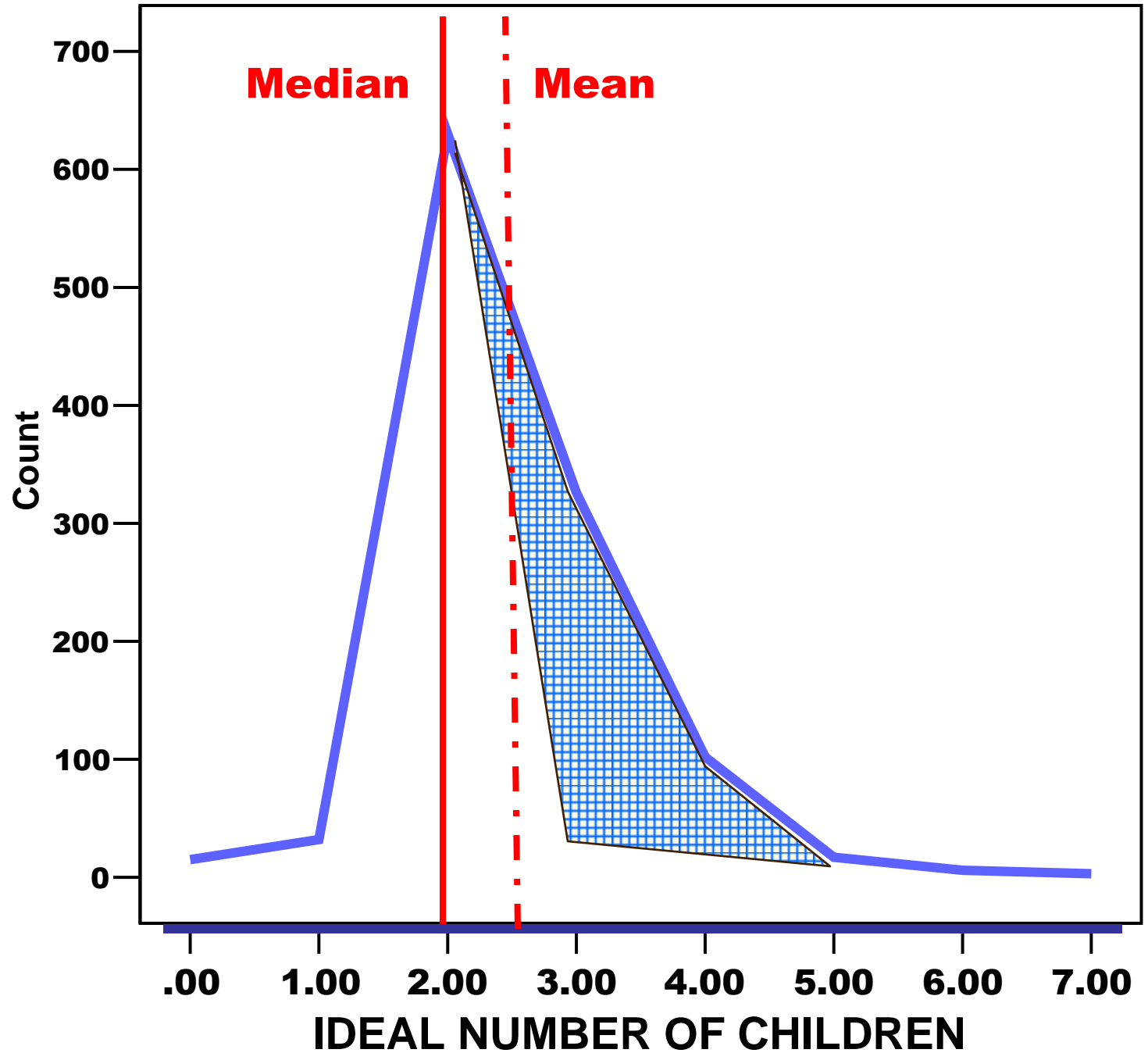- Negative skew has a "tail" to left of Mdn

For most continuous variables, a positively skewed distribution typically has a mean much larger than its median. A negatively skewed distribution typically has a mean smaller than its median.

U.S. household income is positively skewed: in 2006 the median was $48,201 but the mean was $66,570. What produced this gap?

The 2008 GSS asked, "What do you think is the ideal number of children for a family to have?"

Mdn = 2.00     Mean = 2.49     Std dev = 0.88     N = 1,131
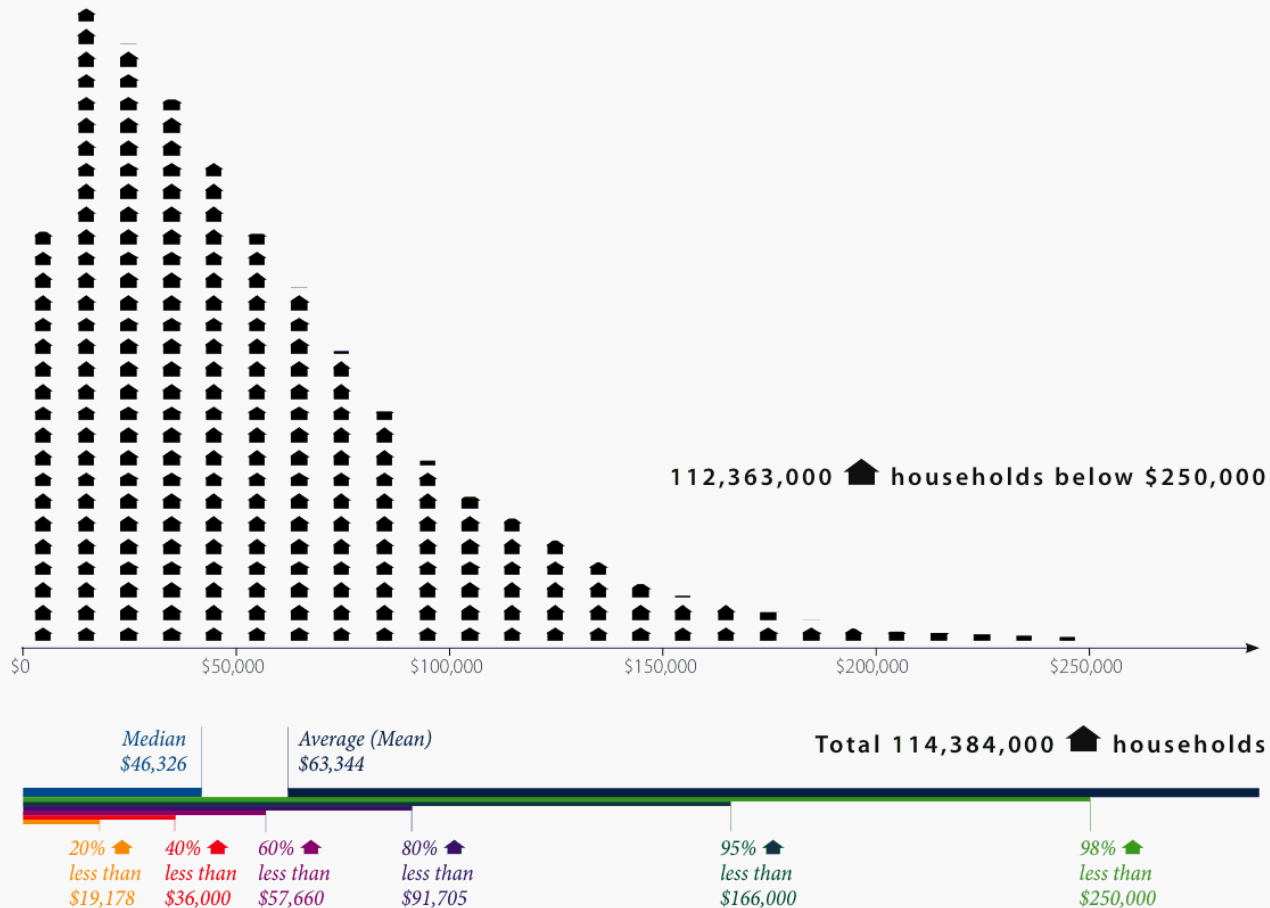
Skewness = _____

# What type of skew does this income distribution have?



Visualizing Economics
Making the *Invisible Hand* Visible

Visit www.visualizingeconomics.com
to view more examples

2005 United States
**Income Distribution (Bottom 98%)**
Each 🏠 equals 500,000 households

112,363,000 🏠 households below $250,000

$0    $50,000    $100,000    $150,000    $200,000    $250,000

Median $46,326    Average (Mean) $63,344

Total 114,384,000 🏠 households

| 20% 🏠 less than $19,178 | 40% 🏠 less than $36,000 | 60% 🏠 less than $57,660 | 80% 🏠 less than $91,705 | 95% 🏠 less than $166,000 | 98% 🏠 less than $250,000 |

# Calculate $s^2$ and $s$ for these 8 ungrouped scores

$$Y_i \quad - \quad \overline{Y} \quad = \quad d_i \qquad (d_i)^2$$

1 - 5 = _____ _____ $\qquad \displaystyle\sum_{i=1}^{8}(d_i)^2 =$ _____

3 - 5 = _____ _____

4 - 5 = _____ _____ $\qquad s_Y^2 = \displaystyle\sum (d_i)^2 \; / (N-1) =$

5 - 5 = _____ _____

6 - 5 = _____ _____

6 - 5 = _____ _____ $\qquad$ _____

7 - 5 = _____ _____

8 - 5 = _____ _____ $\qquad s_Y = \sqrt{s_Y^2} =$ _____

# Calculate variance & standard deviation of NATEDUC

"Are we spending too much money, too little money, or about the right amount on the nation's education system?"

**N = 993**     **Mean = 1.34**

| Category | $Y_i$ | $f_i$ | $(d_i)^2(f_i)$ |
|----------|-------|-------|----------------|
| TOO LITTLE | 1 | 707 | _____ |
| ABOUT RIGHT | 2 | 232 | _____ |
| TOO MUCH | 3 | 54 | _____ |

$$\sum_{i=1}^{K} (d_i)^2(f_i) = \underline{\qquad}$$

$$s_Y^2 = \frac{\sum_{i=1}^{K} (d_i)^2(f_i)}{N-1} = \underline{\qquad\qquad\qquad}$$

$$s_Y = \sqrt{s_Y^2} = \underline{\qquad\qquad}$$

# Calculate variance & standard deviation of SEXFREQ

**N = 1,686**      **Mean = 57.3**

| Category | $Y_i$ | $f_i$ | $(d_i)^2(f_i)$ |
|---|---|---|---|
| NOT AT ALL | 0 | 416 | _____ |
| ONCE OR TWICE | 2 | 149 | _____ |
| ONCE A MONTH | 12 | 176 | _____ |
| 2-3 per MONTH | 36 | 243 | _____ |
| WEEKLY | 52 | 285 | _____ |
| 2-3 per WEEK | 156 | 309 | _____ |
| 3+ per WEEK | 208 | 108 | _____ |

$$\sum_{i=1}^{K} (d_i)^2(f_i) = \underline{\hspace{4cm}}$$

$$s_Y^2 = \frac{\sum_{i=1}^{K} (d_i)^2(f_i)}{N-1} = \underline{\hspace{6cm}}$$

$$s_Y = \sqrt{s_Y^2} = \underline{\hspace{3cm}}$$