# Data: Information or Just a Bunch of Numbers?

## Jim Burati

### Clemson University, Civil Engineering Dept.

*2005 Southeastern Pavement Management & Design Conference*
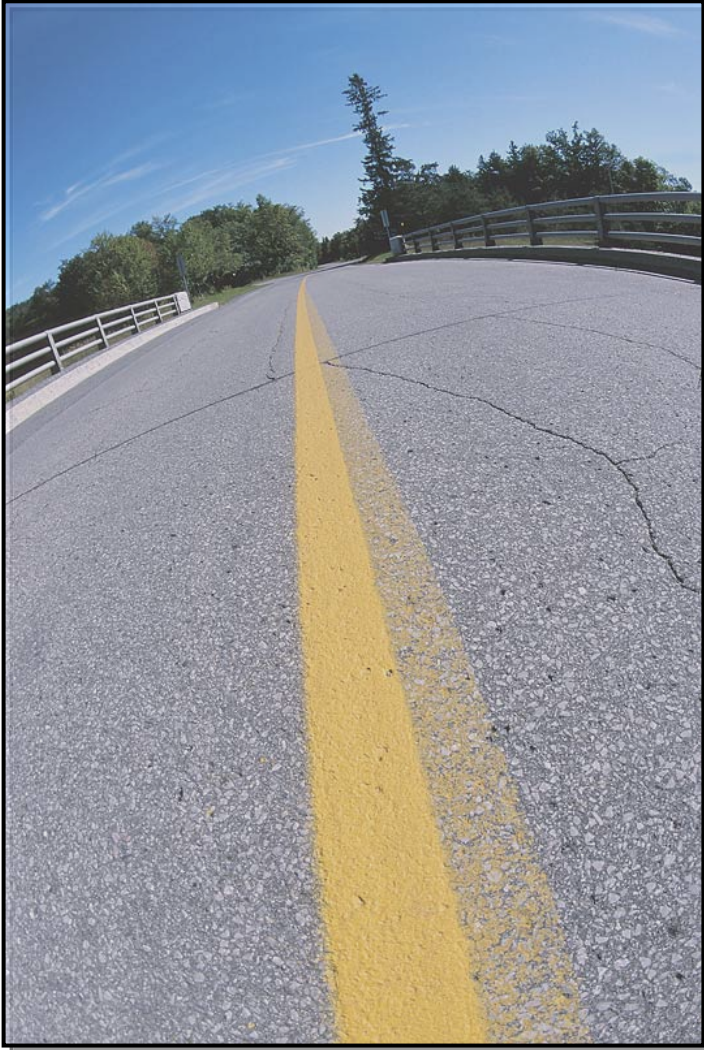
*June 21, 2005*

# Topics

- ❑ Variability
- ❑ Data
- ❑ Estimating Parameters
- ❑ Probability Distributions
- ❑ Drawing Conclusions
- ❑ Regression Analysis?
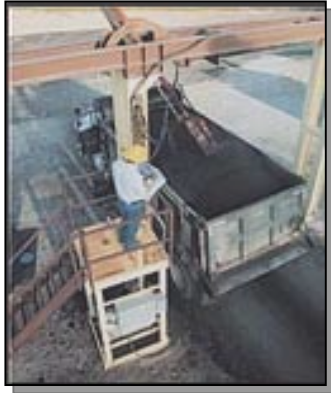
# Sources of Variability

## Material

# Sources of Variability

## Process

# Sources of Variability

*Truck?*

**Sampling**

*Drive?*

*Walk?*

*Total Project?*

*Portion?*

*Road?*

# Sources of Variability

## Testing

# **Variability**

❑ Which variability do we obtain? Which do we need?

➢ **Material variability.** $\left.\vphantom{\begin{matrix} a \\ b \end{matrix}}\right\}$ $\sigma_m^2$
➢ **Process variability.**

➢ **Sampling variability.** $\sigma_s^2$

➢ **Testing variability.** $\sigma_e^2$

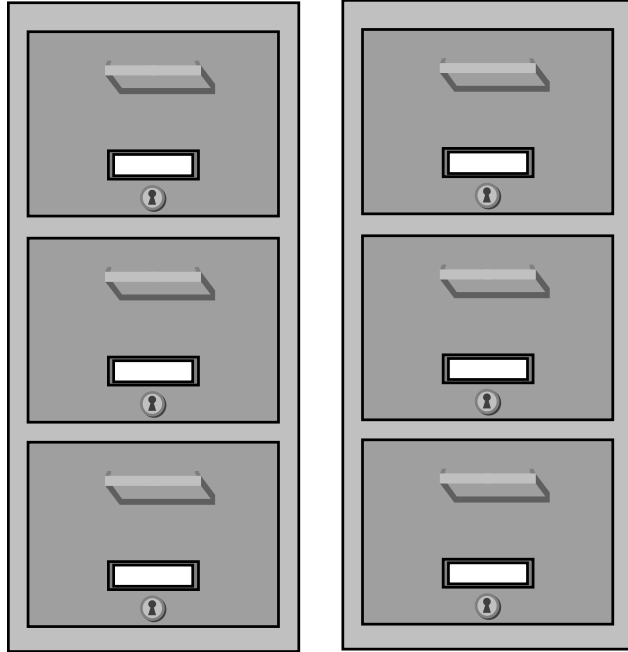# Overall Variability

*Obtain?*

$$\sigma_o^2 = \sigma_m^2 + \sigma_s^2 + \sigma_e^2$$

*Need?*

$$\sigma_o^2 = \boxed{\sigma_m^2} + \sigma_s^2 + \sigma_e^2$$

# Sources of Data

**Be Wary of Historical Records**

*Sampling?*

*Biased Reporting?*

*Test Methods?*

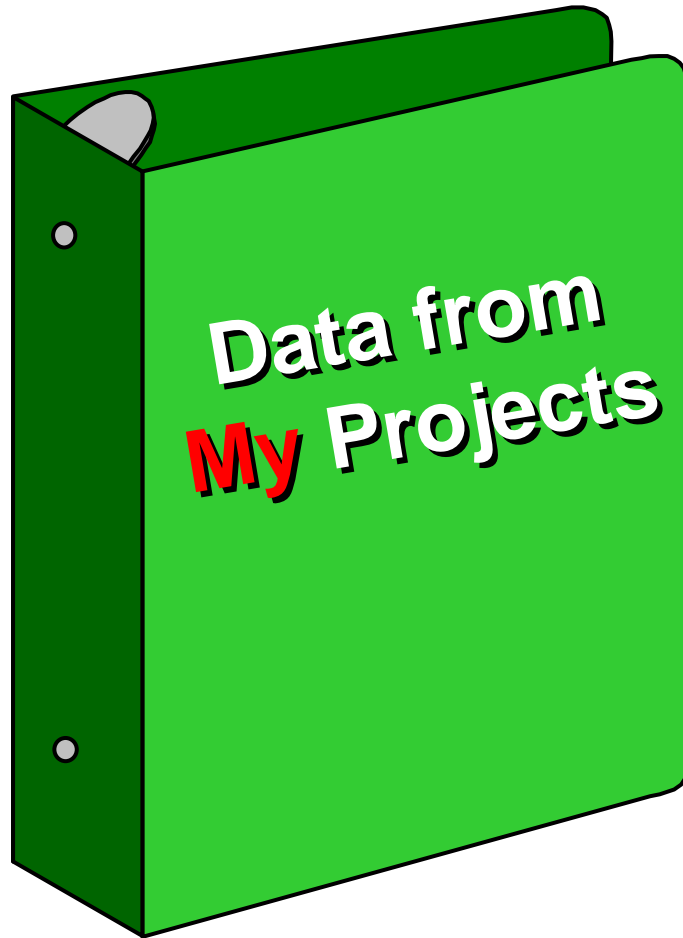# Sources of Data

**Data from Other States?**

**FHWA?**

*Materials?*

*Procedures?*
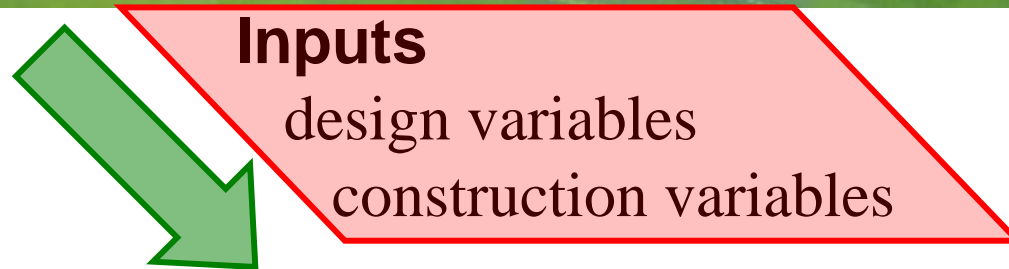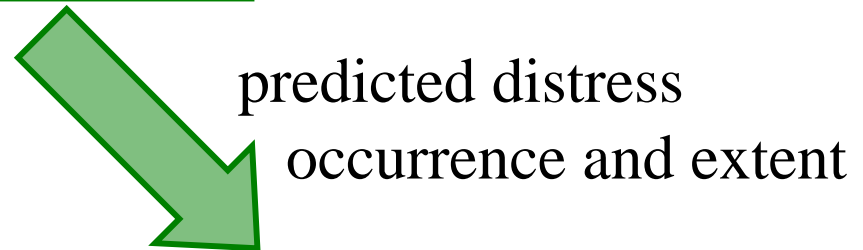
DOT Research

**Data from My Projects**

**Random Sampling
&
Unbiased Reporting
&
Same Procedures**

# Performance Related Specs

**Inputs**

design variables

construction variables

**Performance-Prediction Models**

predicted distress

occurrence and extent

**Maintenance-Cost Models**

**Peter Kopac**
**"Making Roads Better and Better"**
***Public Roads,* Vol. 66, No. 1, 2002**

**Outputs**

life-cycle costs

# Mechanistic Empirical Model

**Inputs**

**Mechanistic-Empirical Model – 1, 2a, 2b, 3**

# Pavement Design

# What do we want to know?

→ **Center** ←

← **Spread** →

# What Data Do We Need?

- ❑ **Parameters**
  - ➢ **Mean.**
  - ➢ **Standard Deviation.**
- ❑ **Probability Distribution**

# Things are easy, right?

- ❑ **Parameters**
  - ➢ **Center? Mean.** $\bar{x}\ for\ \mu$
  - ➢ **Spread? Standard Deviation** $s\ for\ \sigma$

- ❑ **Probability Distributions**
  $Normal$

# Example

- ❑ Normal Population with:
  - – **Mean = 100**
  - – **Standard Deviation = 10**
- ❑ Wish to estimate the mean and standard deviation.
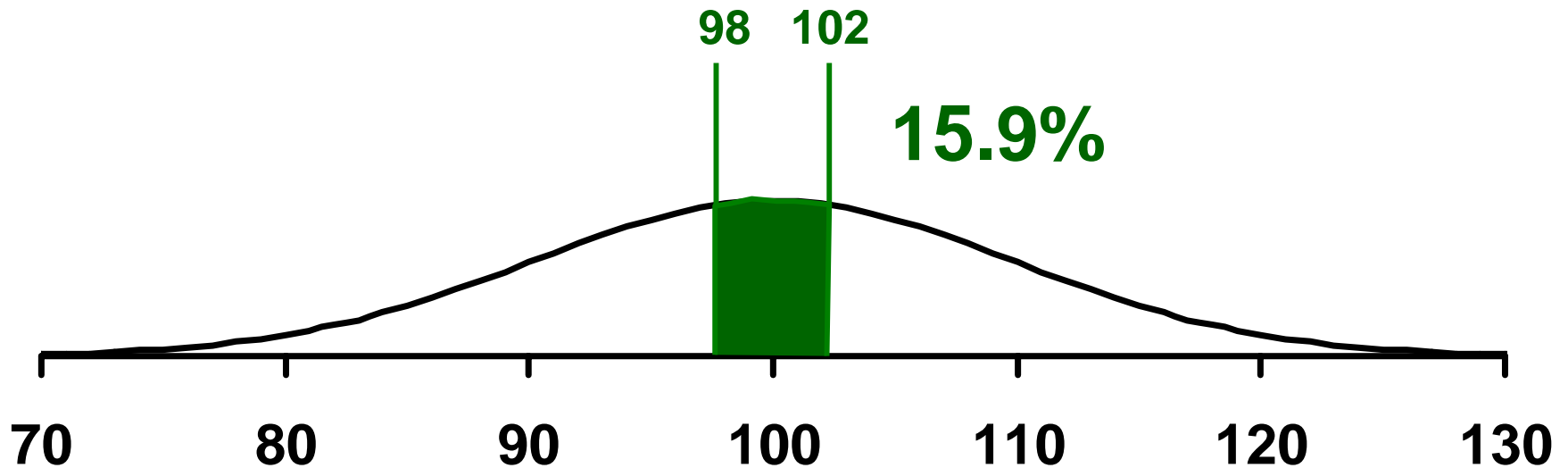
# Estimating the Mean

❑ **Unbiased estimate**

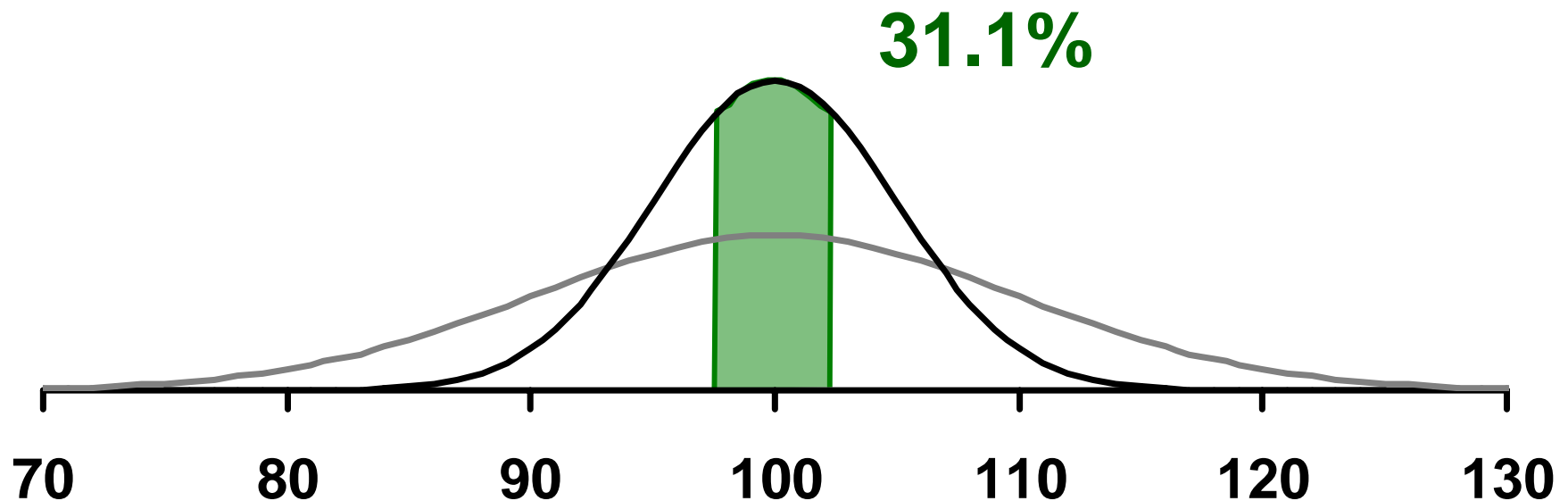$$\overline{X} \text{ for } \mu$$

❑ **Confidence Interval**

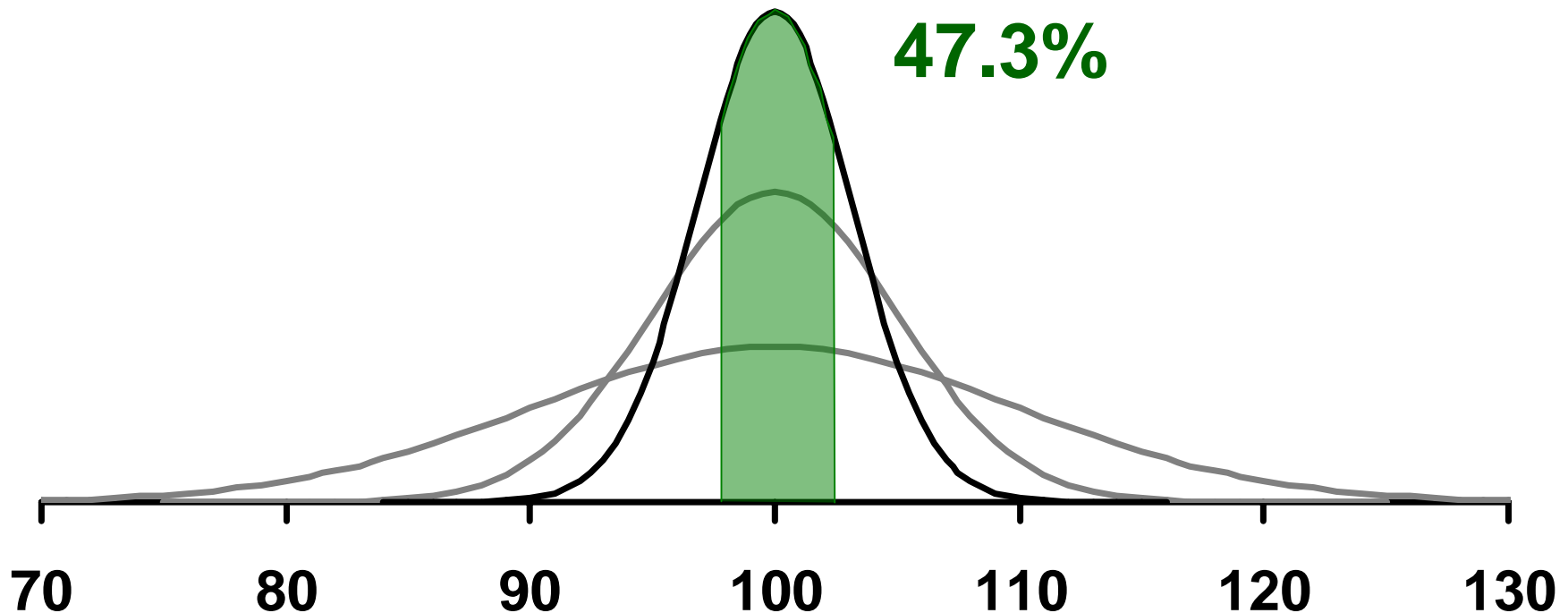$$Sample\ Size$$

# Central Limit Theorem, $n = 1$

$\sigma = 10$



**98  102**

**15.9%**

70    80    90    100    110    120    130

# Central Limit Theorem, $n = 4$

$$\sigma_{\overline{X}} = \frac{\sigma}{\sqrt{n}} = \frac{\sigma}{\sqrt{4}} = 5$$

**31.1%**

70　　　80　　　90　　　100　　　110　　　120　　　130

# Central Limit Theorem, $n = 10$

$$\sigma_{\overline{X}} = \frac{10}{\sqrt{10}} \cong 3.16$$

**47.3%**

70        80        90        100        110        120        130

# Central Limit Theorem, $n = 20$

$$\sigma_{\overline{X}} = \frac{10}{\sqrt{20}} \approx 2.24$$

**62.9%**

70    80    90    100    110    120    130

# Estimating the Standard Deviation

❑ **Unbiased estimate**

$$s^2 \text{ for } \sigma^2 \qquad s \text{ for } \sigma \; ?$$

❑ **Confidence Interval**

$$Sample\ Size$$

# Estimator for $\sigma$ ?

$$s^2 = \frac{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}{n-1}$$

$$s = \sqrt{\frac{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}{n-1}}$$

# Sampling Distribution, $s^2$

$$\frac{(n-1)s^2}{\sigma^2} \Rightarrow \chi^2 - \text{distribution}$$
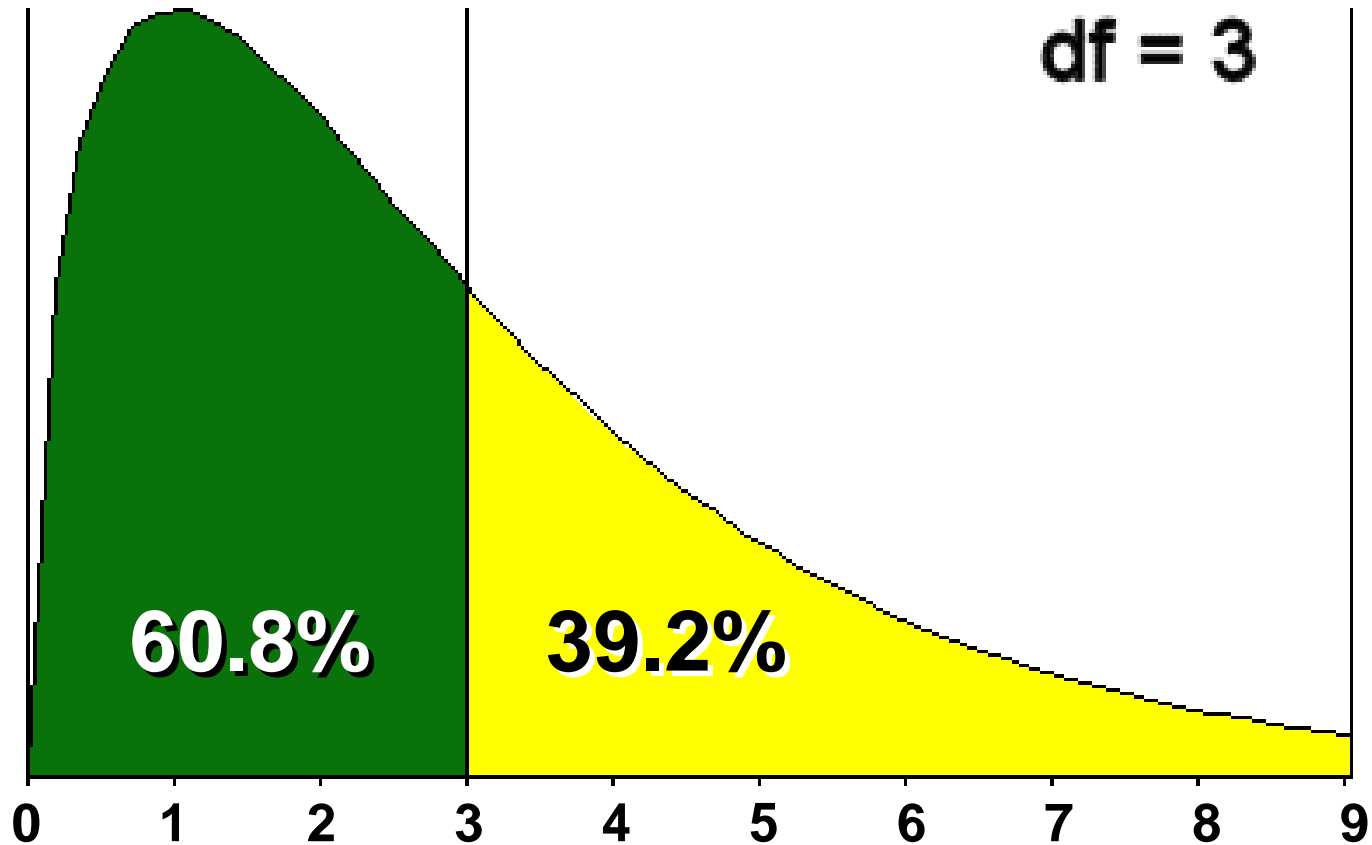
Since $s^2$ is an unbiased estimator for $\sigma^2$,

The average of $s^2/\sigma^2$ equals 1, so the mean of the $\chi^2$ distribution $= (n-1)$.
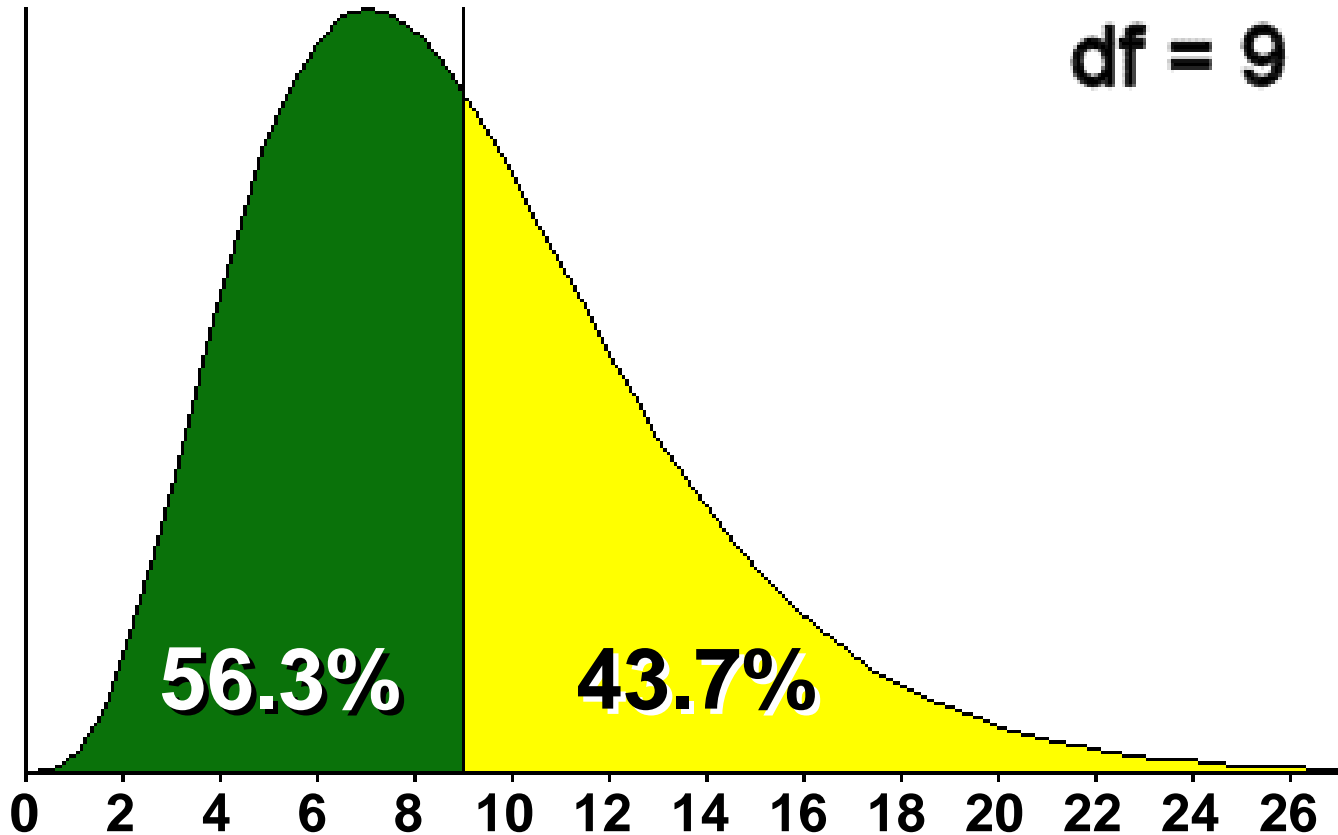
# $\chi^2$ Distribution

❑ Shape varies with degrees of freedom, i.e., $n - 1$.

❑ It is always positive.
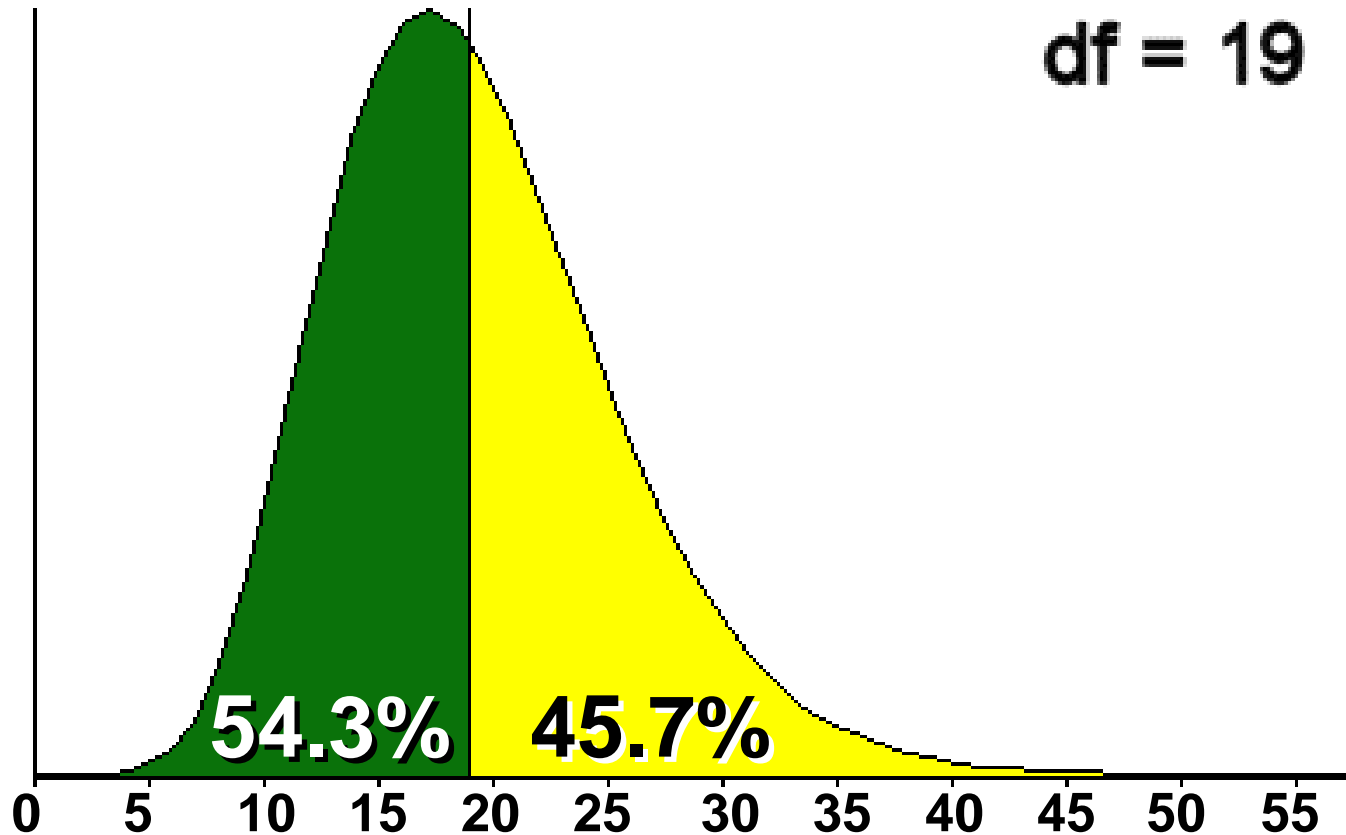
❑ It is never symmetric.

# $\chi^2$ Distribution, $n = 4$



df = 3

60.8%    39.2%

0   1   2   3   4   5   6   7   8   9

# $\chi^2$ Distribution, $n = 10$

df = 9

56.3%  43.7%

0  2  4  6  8  10  12  14  16  18  20  22  24  26

# $\chi^2$ Distribution, $n = 20$



df = 19

54.3%    45.7%

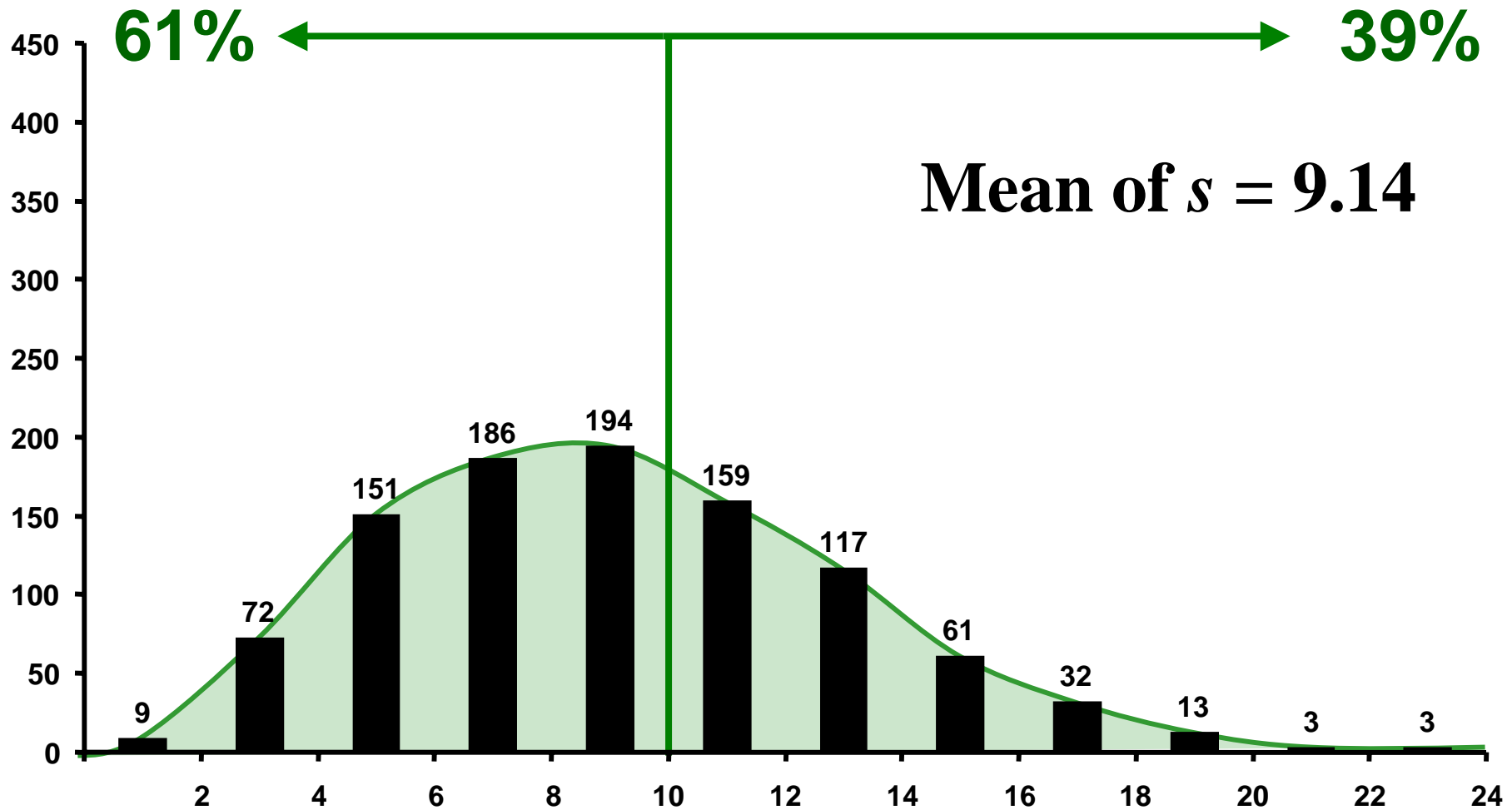0  5  10  15  20  25  30  35  40  45  50  55

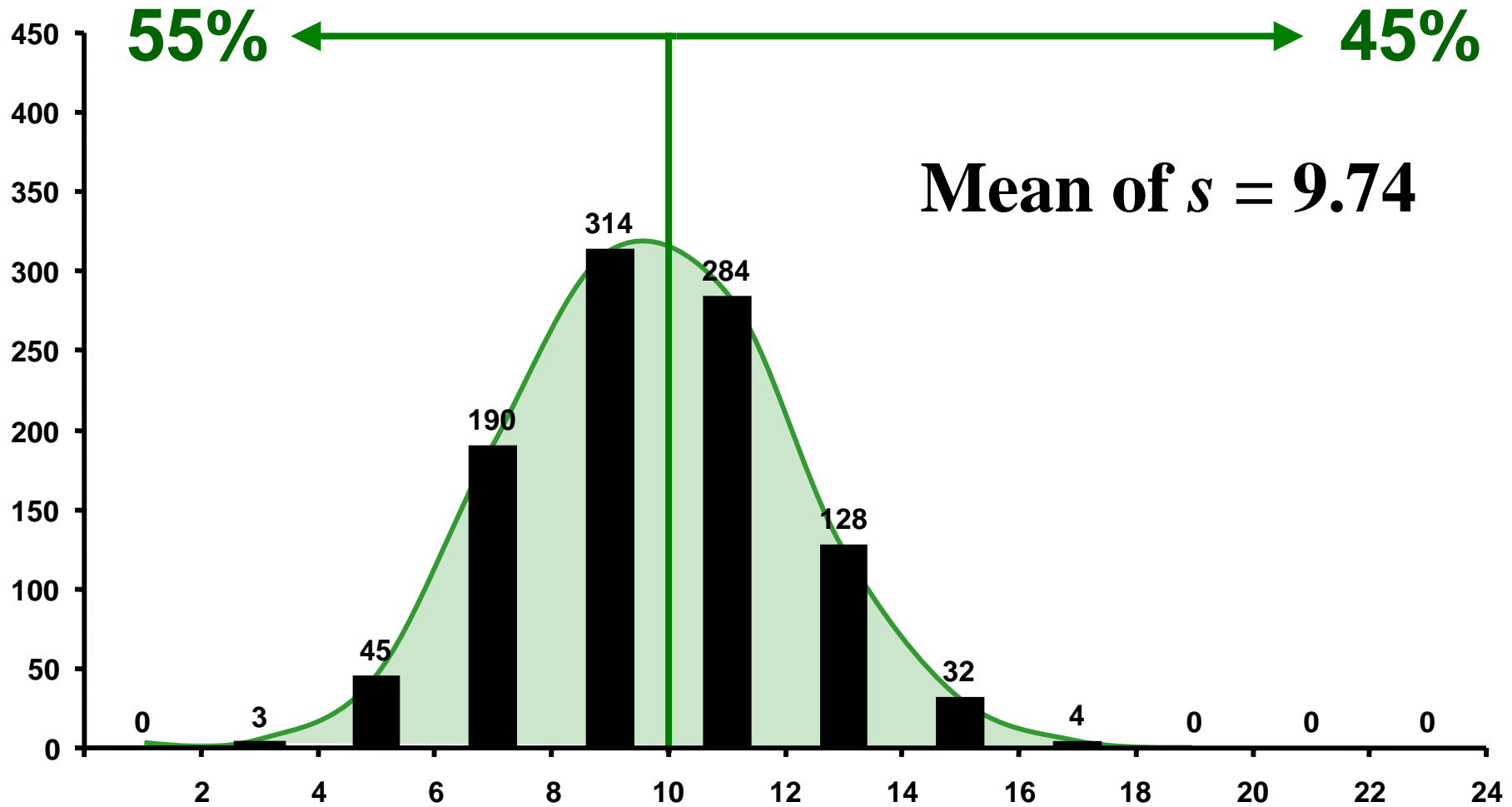*Note,* for df = 200:    51.3%  48.7%

# Sampling Example

❑ Normal Population with:

   – **Mean = 100**

   – **Std. Dev. = 10 (Variance = 100)**

❑ Draw 1000 Samples

   – $n = 4, 10, 20$

   – **Calculate each** $s$
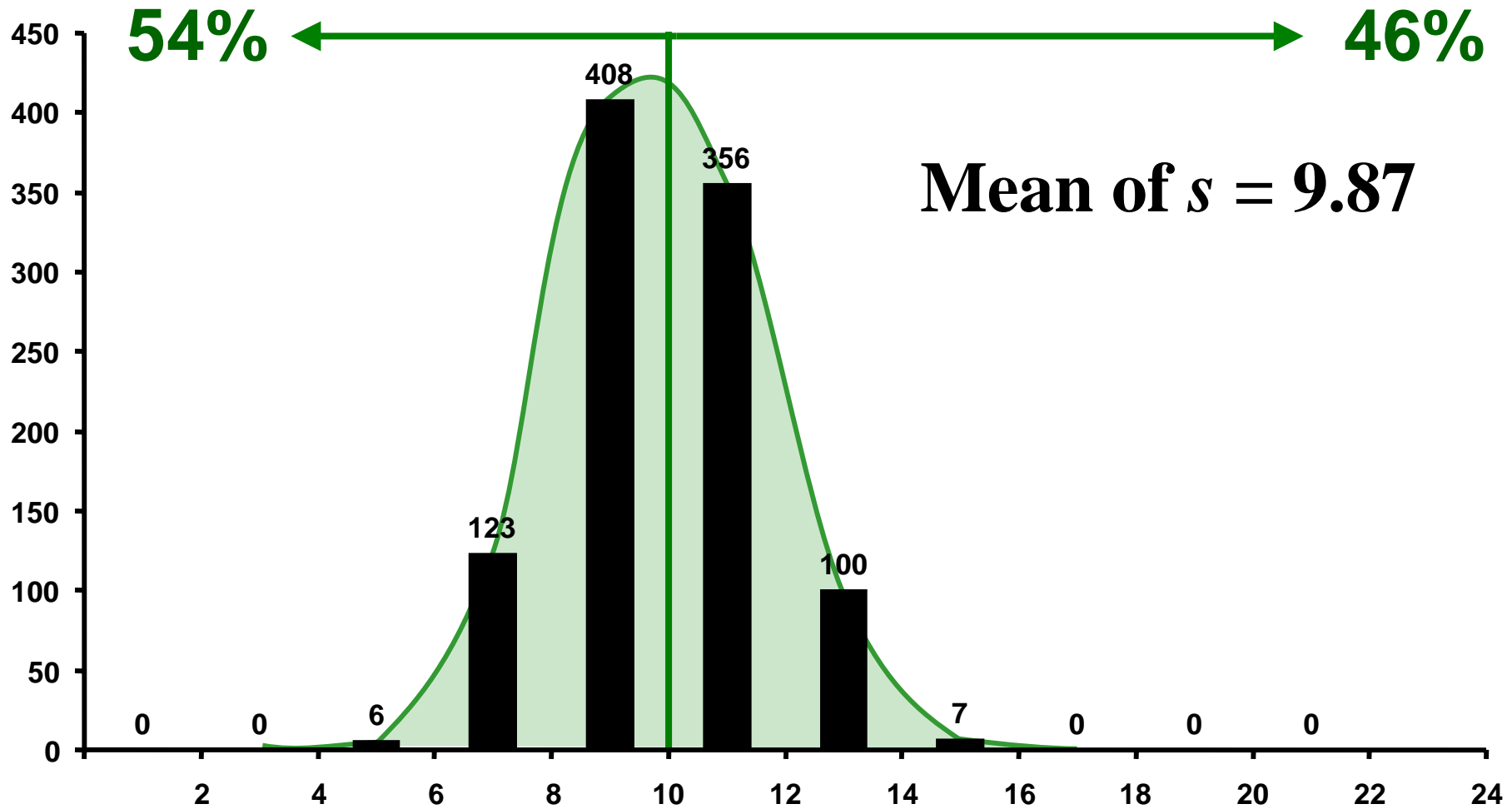
# Estimating $\sigma$ with $s$, $n = 4$



**61%** ← → **39%**

Mean of $s = 9.14$

# Estimating $\sigma$ with $s$, $n = 10$



**55%** ← → **45%**

Mean of $s = 9.74$

Histogram values (left to right): 0, 3, 45, 190, 314, 284, 128, 32, 4, 0, 0, 0

# Estimating $\sigma$ with $s$, $n = 20$



**54%** ← → **46%**

Mean of $s = 9.87$

Histogram values: 0, 0, 6, 123, 408, 356, 100, 7, 0, 0, 0

X-axis: 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24

Y-axis: 0, 50, 100, 150, 200, 250, 300, 350, 400, 450

# Estimators for $\sigma$ ?

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

**Biased Low!**

# Bias < 1% when $n > 26$

# Distribution of $s$, using $N$ *d.f.*

$$s \equiv \sqrt{\frac{1}{N}\sum_{i=1}^{N}(x_i - \overline{x})^2} \qquad \sigma^2 \equiv \frac{Ns^2}{N-1}$$

$$f_N(s) = 2\frac{\left(\dfrac{N}{2\sigma^2}\right)^{(N-1)/2}}{\Gamma\left(\dfrac{1}{2}(N-1)\right)} e^{-Ns^2/(2\sigma^2)} s^{N-2}$$

# Distribution of *s*, using N *d.f.*



$f_N(s)$

2

4

8

12

*S*

# Unbiased Estimator for $\sigma$

$$(c(n) \text{ or } c_4)s \left[ \frac{\Gamma\left(\frac{n-1}{2}\right)\sqrt{n-1}}{\Gamma\left(\frac{n}{2}\right)\sqrt{2}} \right]s$$

**Wow!**

$$\text{where} \quad \Gamma(z) \equiv \int_0^\infty t^{z-1} e^{-t} dt$$

# Unbiased Estimator for $\sigma$

$$c(n) = \begin{cases} \dfrac{2^{n-2.5}\sqrt{n-1}}{(n-2)\sqrt{\pi}\dbinom{n-3}{\tfrac{n-3}{2}}} & n \text{ odd} \\[2em] \dbinom{n-3}{\tfrac{n-2}{2}}\dfrac{\sqrt{\pi(n-1)}}{2^{n-2.5}} & n \text{ even}, n \geq 4 \end{cases}$$

**Wow!**

# Unbiased Estimator for $\sigma$

$$c's \qquad \left( \frac{n - 0.75}{n - 1} \right) s$$

**Accurate to ¼ %**
**for all $n$.**

# Sampling Example

- ❑ Normal Population with:
  - **Mean = 100**
  - **Std. Dev. = 10 (Variance = 100)**
- ❑ Draw 1000 Samples
  - $n = 4, 10, 20$
  - **Calculate each** $c's$

# Estimating $\sigma$ with $c's$, $n = 4$

**53%** ⟵————————————⟶ **47%**

Mean of $c's = 9.92$

Mean of $s = 9.14$

Histogram data (value : frequency):
- 1 : 7
- 3 : 59
- 5 : 126
- 7 : 167
- 9 : 171
- 11 : 178
- 13 : 120
- 15 : 80
- 17 : 50
- 19 : 28
- 21 : 9
- 23 : 2
- 25 : 3

Y-axis: 0, 50, 100, 150, 200, 250, 300, 350, 400, 450

X-axis: 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24, 26

# Estimating $\sigma$ with $c's$, $n = 10$



**50%**     ← ———————————— →     **50%**

Mean of $c's = 10.02$

Mean of $s = 9.74$

# Estimating $\sigma$ with $c$'s, $n = 20$

**50%** ← → **50%**

Mean of $c$'s = 10.00

Mean of $s$ = 9.87

Histogram values (count by x-value):
- 0 (at 2)
- 0 (at 4)
- 5 (at ~5)
- 112 (at 7)
- 386 (at 9)
- 381 (at 11)
- 105 (at 13)
- 11 (at 15)
- 0 (at 18)
- 0 (at 20)
- 0 (at 22)
- 0 (at 24)
- 0 (at 26)

Y-axis: 0, 50, 100, 150, 200, 250, 300, 350, 400, 450
X-axis: 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24, 26

❑ **Unbiased estimate**

$$c's \quad versus \quad s$$

# Probability Distribution

❑ Normality Assumption?

❑ Approximately Normal?

❑ Skewed? *(median)*

❑ Goodness of Fit to Normal?

– **Histogram.**

– **Normal probability paper.**

– **Statistical test: ~~Chi Square, K-S?~~**

– **Shapiro-Wilk, Anderson-Darling**

# Skewed Distributions

**Positive**

**Negative**

# Pavement Materials

## Few materials have skewed distributions.

**0**

**However**

# Outliers and Skewness

**ASTM E-178**

# Bimodal Distributions

# How?

# Obvious

# Not So Obvious?

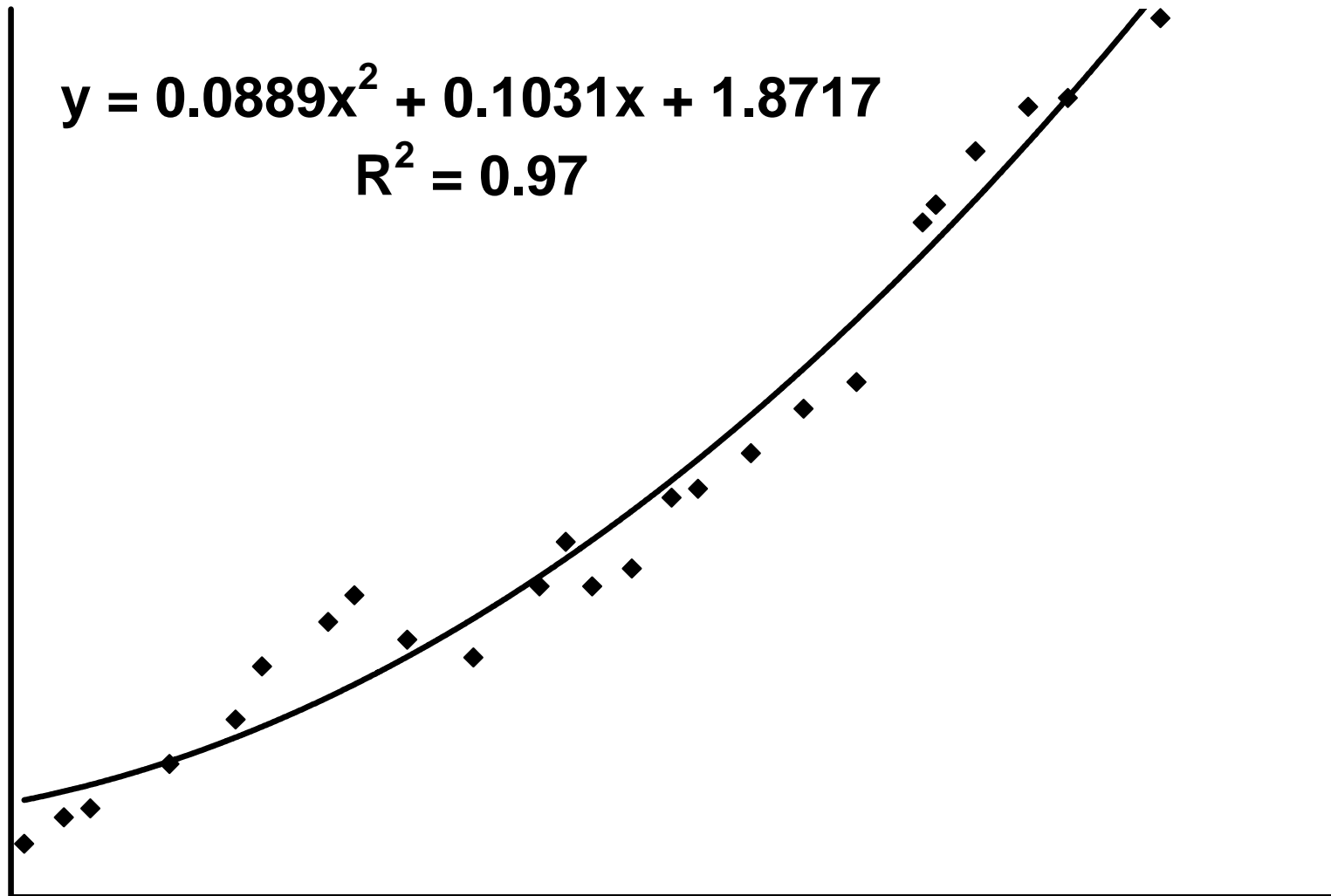# Multiple Raters?

# Drawing Conclusions from Data?

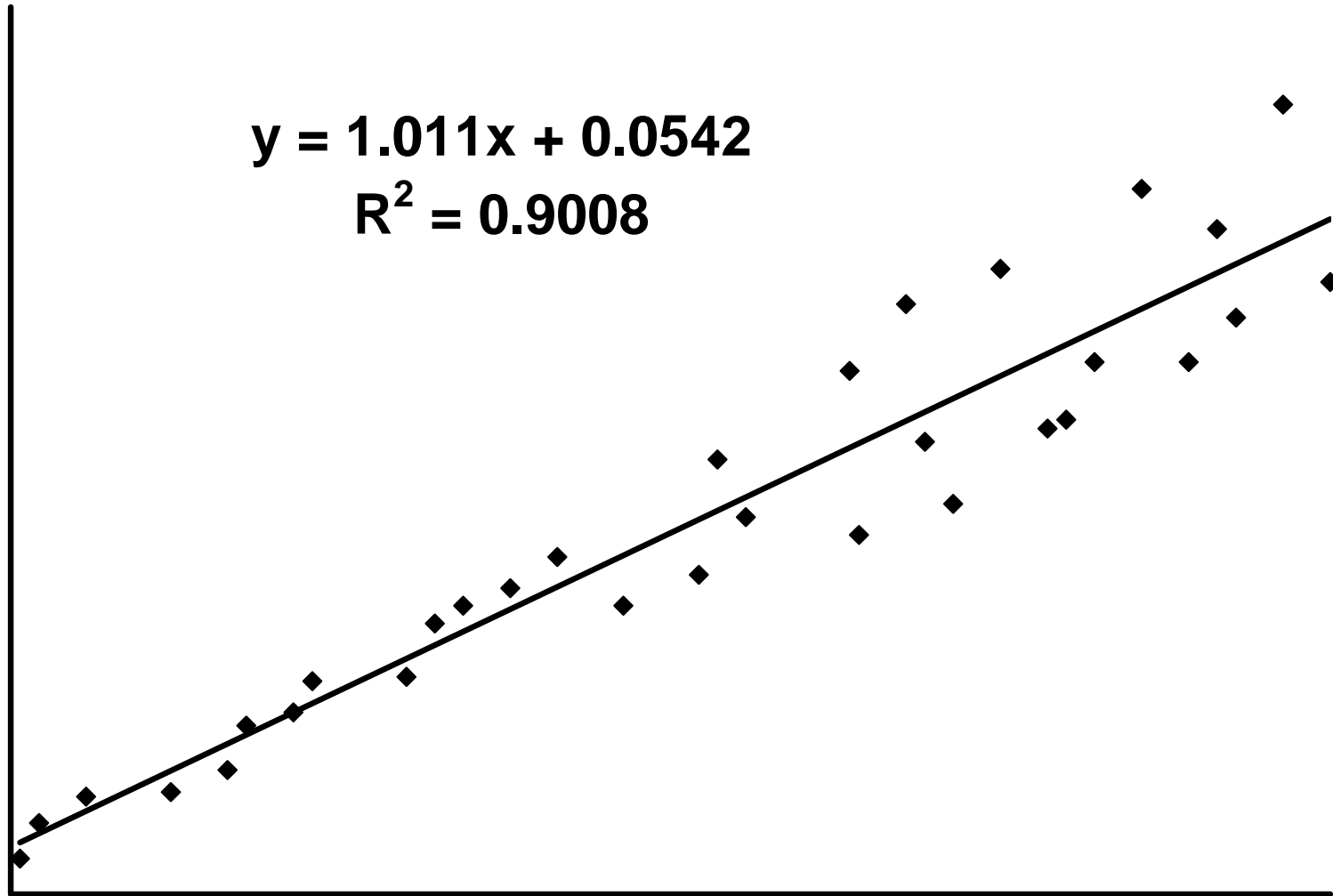# Regression Lines?

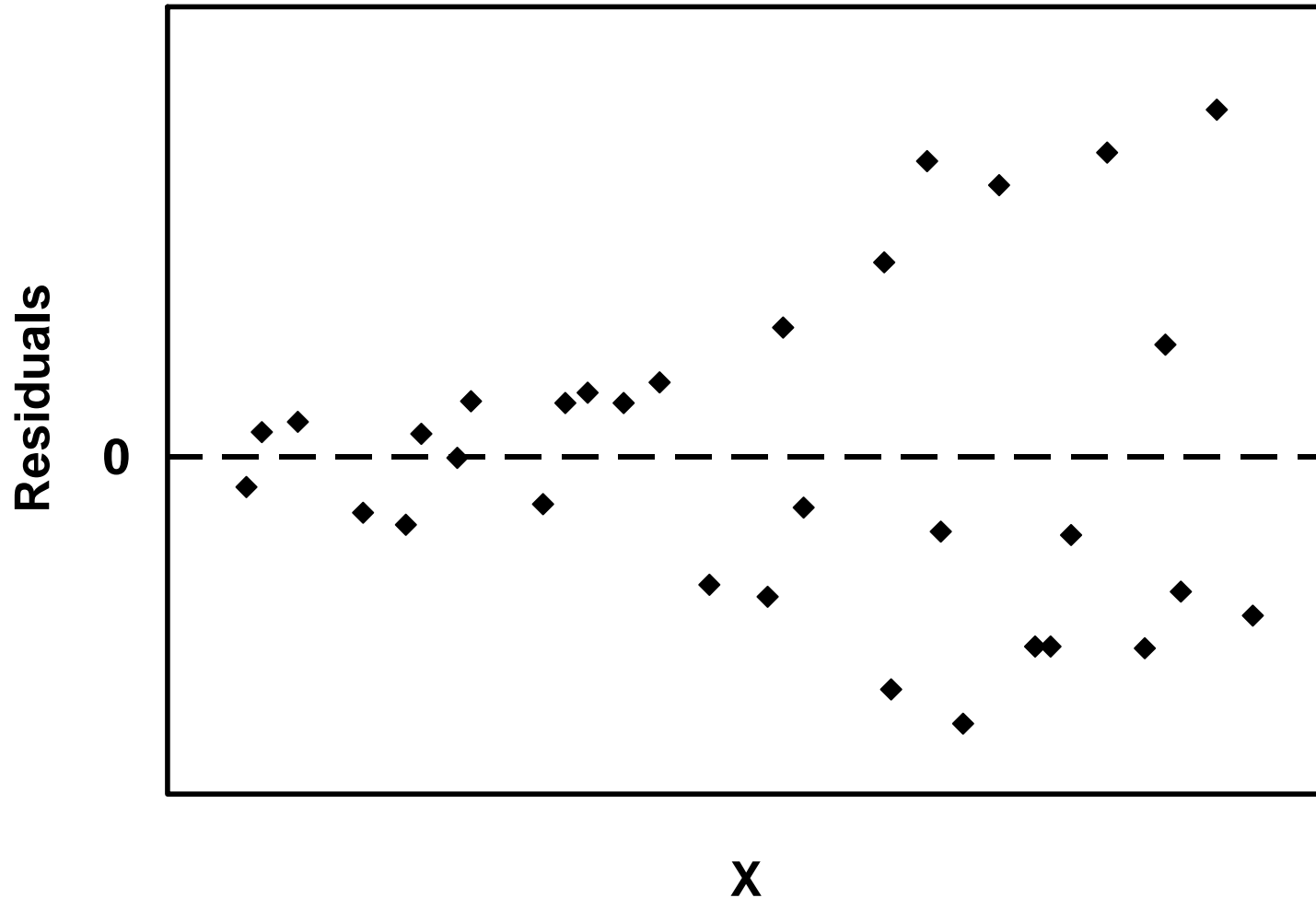**Assumptions?**

# Regression & Normality?

# Residuals Normal?



Residuals

0.0

X

# Correct Model?

$$y = 0.0889x^2 + 0.1031x + 1.8717$$
$$R^2 = 0.97$$

# Regression & Variability?



$$y = 1.011x + 0.0542$$
$$R^2 = 0.9008$$

# Regression & Variability?

# Classic Example

**For Each of Four $(x, y)$ Data Sets:**

$$n = 11 \qquad\qquad r = 0.82$$

$$\bar{x} = 9.0 \qquad\qquad y = 0.50x + 3.00$$

$$\bar{y} = 7.5 \qquad\qquad R^2 = 0.67$$

$$s_x = 3.32$$

$$s_y = 2.03$$

# Data Set 1


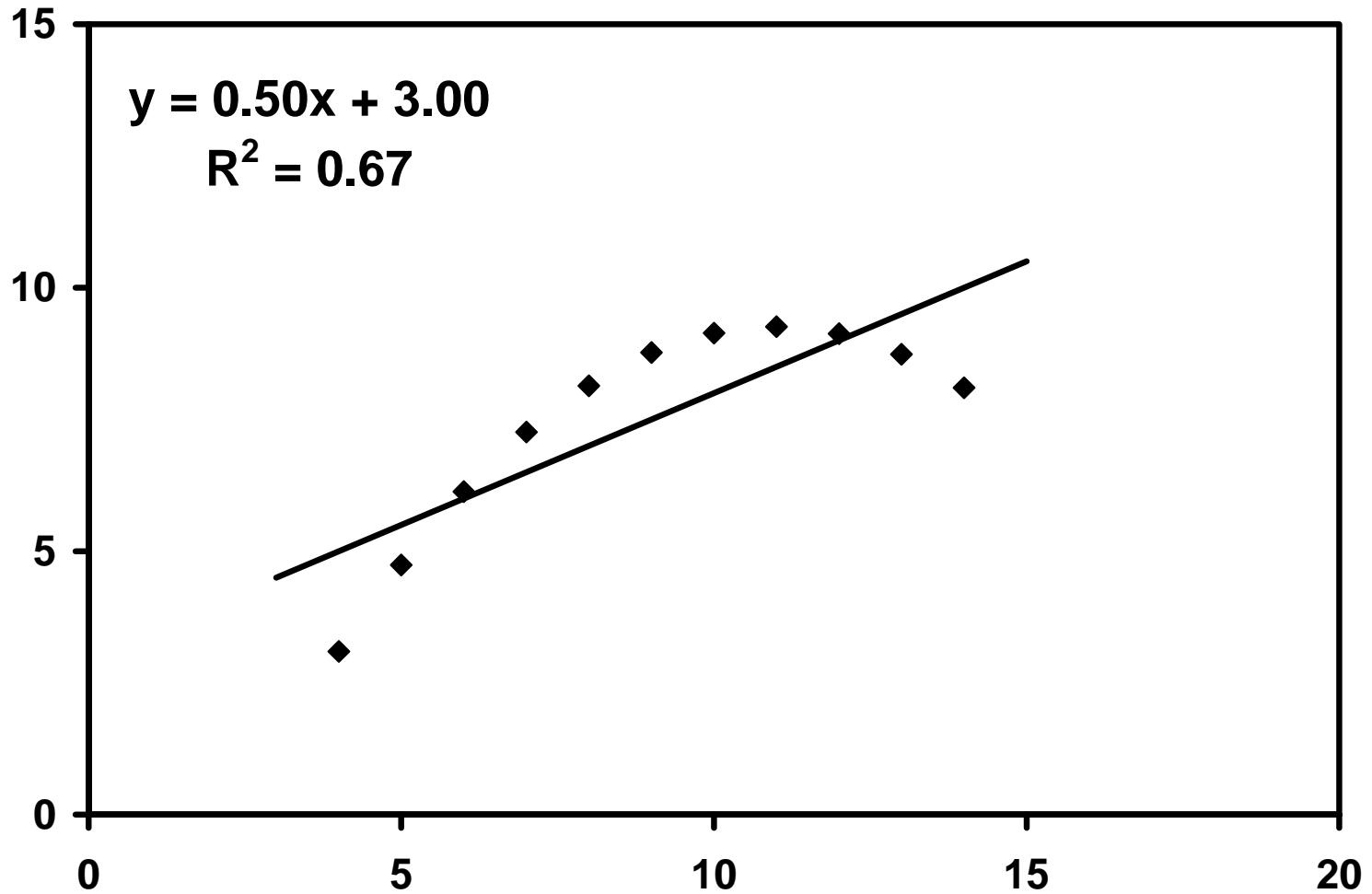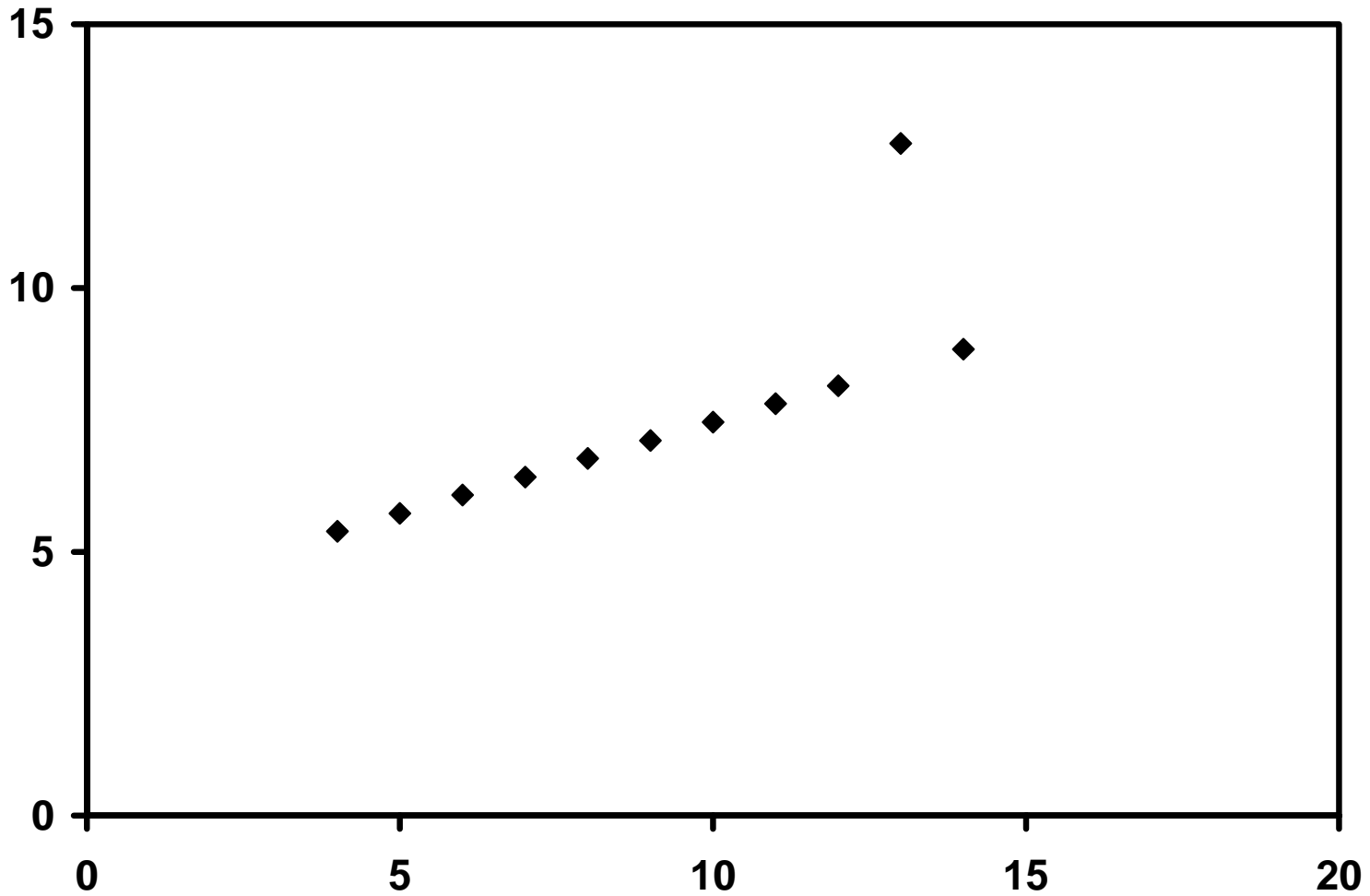
$y = 0.50x + 3.00$

$R^2 = 0.67$

# Data Set 2

# Data Set 2



$y = 0.50x + 3.00$

$R^2 = 0.67$

# Data Set 3

$$y = 0.50x + 3.00$$
$$R^2 = 0.67$$

# Data Set 4

# Data Set 4

$$y = 0.50x + 3.00$$
$$R^2 = 0.67$$
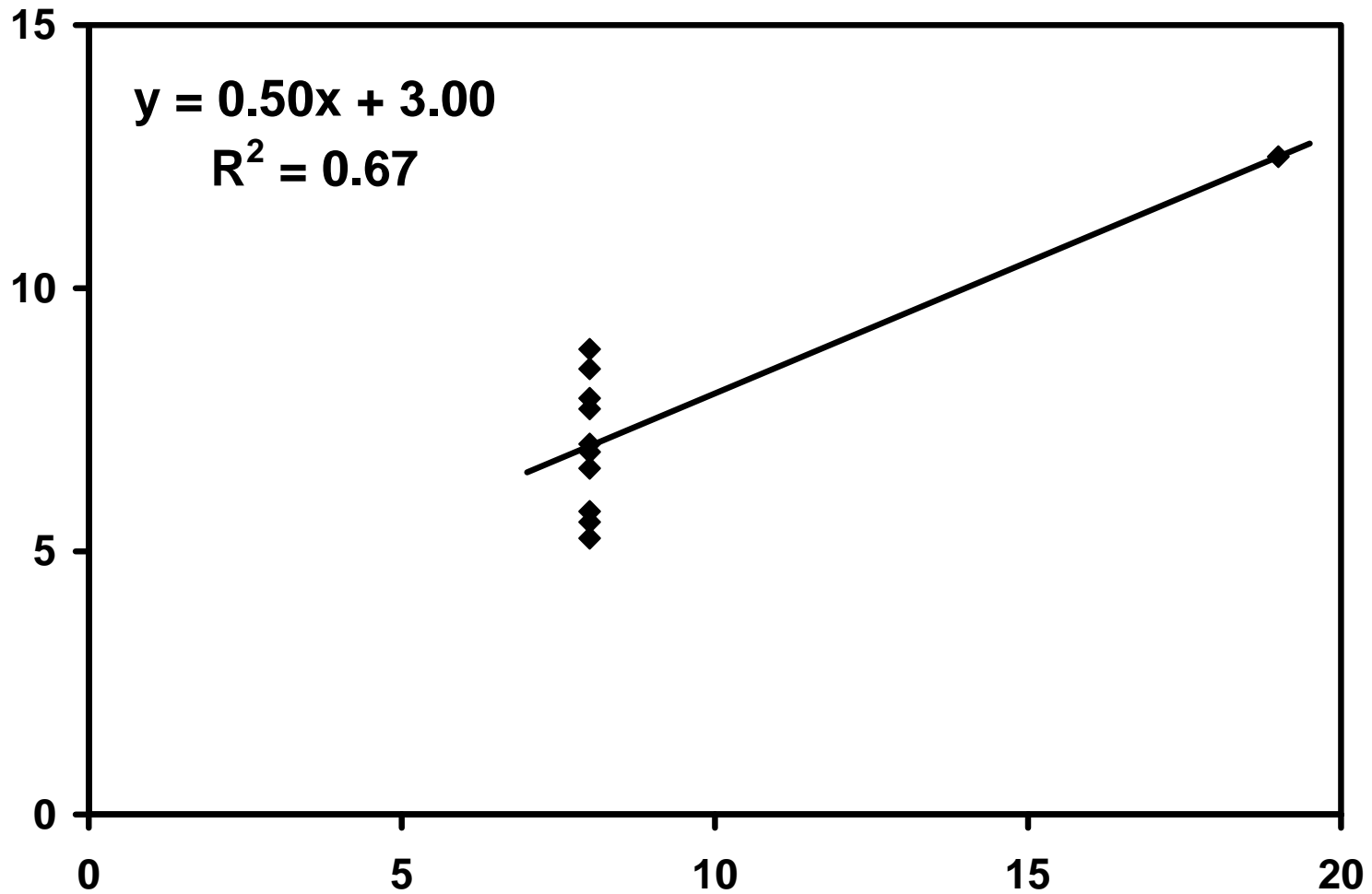
# Closing

❑ Getting data is easy.

❑ Getting valid data is not as easy.

❑ Analyses limited by quality of data.

❑ Implicit assumptions (e.g., normal).

❑ Data

– Can produce meaningful decisions
– Can be meaningless numbers
– Can lead to erroneous conclusions.

# *Thanks for having me!*

**The End**

**2005 SE Pavement Management & Design Conference**