

# Analytical chemistry

**Statistical data treatment and evaluation. Random errors in chemical analysis**

*"Facts are stubborn, but statistics are much more pliable."*

—Mark Twain

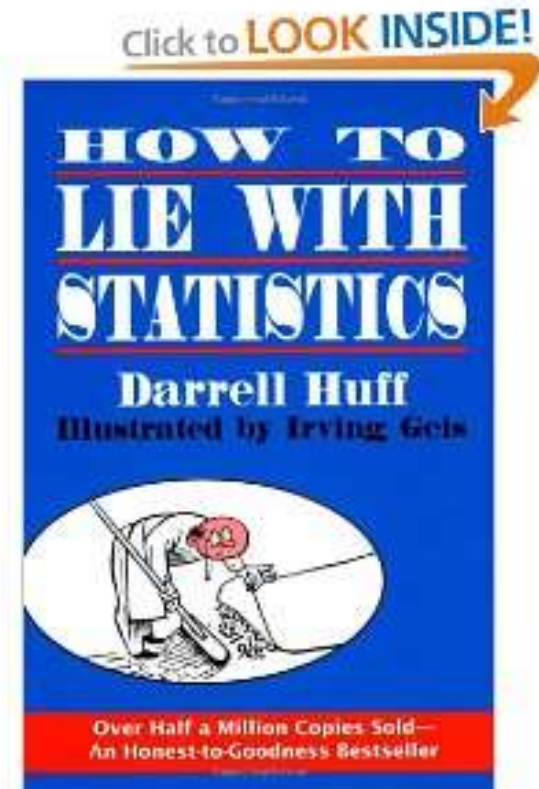
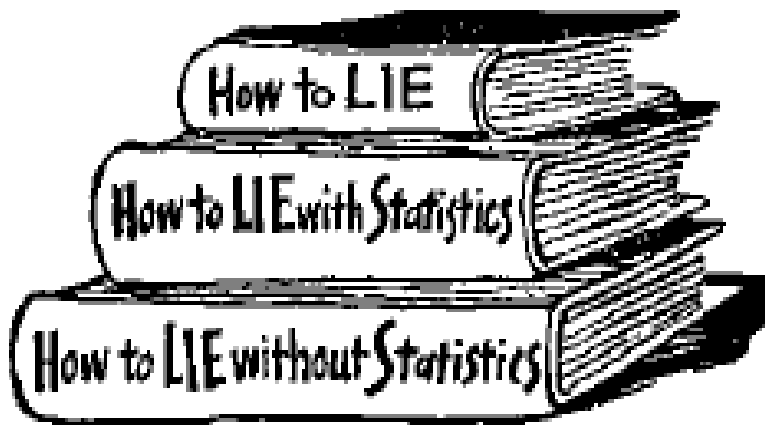


**No  
1  
1**

**L  
e  
c  
t  
u  
r  
e**

# Statistics - Lying without sinning?

- "Lies, damned lies, and statistics"



1954



Although data handling normally follows the collection of analytical data. A knowledge of statistical analysis will be required as you perform experiments in the laboratory. Also, statistical analysis is necessary to understand the significance of the data that are collected and thus sets limits on each step of the analysis. Experimental design (including required sample size, measurement accuracy, and number of analyses needed) relies on a proper understanding of what the data represent.

Statistical analysis only reveals information that is already present in a data set. *No new information is created by* statistical treatments. Statistical methods, do allow us to categorize and characterize data in different ways and to make objective and intelligent decisions about data quality and interpretation.

It is impossible to perform a chemical analysis that is totally free of errors or uncertainties. We can only hope to minimize errors and estimate their size with acceptable accuracy.



# Significant Figures

**The significant figures** in a number are all of the certain digits plus the first uncertain digit. The digits of a number which are needed to express the precision of the measurement from which the number was derived are known as significant figures

A zero may or may not be significant depending on its location in a number. A zero that is surrounded by other digits is always significant (such as in 30.24 mL) because it is read directly and with certainty from a scale or instrument readout. On the other hand, zeros that only locate the decimal point for us are not. If we write 30.24 mL as 0.03024 L, the number of significant figures is the same. The only function of the zero before the 3 is to locate the decimal point, so it is not significant. Terminal, or final, zeros may or may not be significant. For example, if the volume of a beaker is expressed as 2.0 L, the presence of the zero tells us that the volume is known to a few tenths of a liter so that both the 2 and the zero are significant figures. The number of significant figures in a measurement is independent of the placement of the decimal point. Take the number 92,067. This number has five significant figures, regardless of where the decimal point is placed. For example, *92,067  $\mu\text{m}$ , 9.2067 cm, 0.92067 dm, and 0.092067m* all have the same number of significant figures.





List the proper number of significant figures in the following numbers and indicate which zeros are significant.

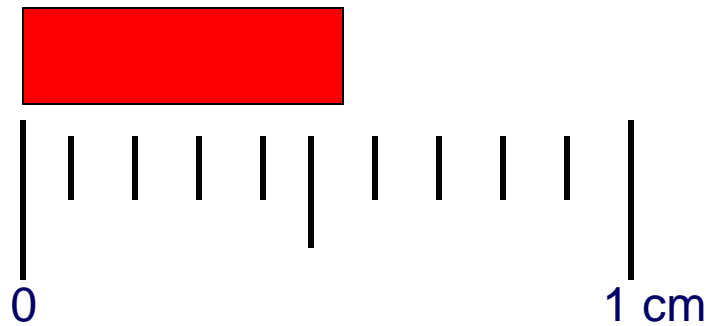
0.216; 90.7; 800.0; 0.0670

### Rules for determining which digits are significant:

1. All non-zero numbers are significant.
2. Zeros between non-zero numbers are significant.
3. Zeros to the right of the non-zero number **and** to the right of the decimal point are significant.
4. Zeros before non-zero numbers are **not** significant.



When reading the scale of any apparatus, you should interpolate between the markings. It is usually possible to estimate to the nearest tenth of the distance between two marks.



**0.55 cm?**

**0.56 cm?**

**2 signif. figs**

**0.55 cm implies an error of at least  $0.55 \pm 0.01$  cm**

**In experimental data, the first uncertain figure is the last significant figure.**



## Significant Figures in Numerical Computations

For *addition and subtraction*, the number of significant figures can be found by visual inspection. For example, in the expression

$$3.4 + 0.020 + 7.31 = 10.730 \text{ (round to } 10.7) = 10.7 \text{ (rounded)}$$

the second and third decimal places in the answer cannot be significant because 3.4 is uncertain in the first decimal place. Hence, the result should be rounded to 10.7.

We can generalize and say that, for *addition and subtraction*, the result should have **the same number of decimal places** as the number with the *smallest number of decimal places*.

- *Logarithms and Antilogarithms*

$$\text{pH} = -\log 2.0 \times 10^{-3} = -(-3 + 0.30) = 2.70$$

$$\log 4.000 \times 10^{-5} = -4.3979 \quad \text{antilog } 12.5 = 3 \times 10^{12}$$

1. In a logarithm of a number, keep as many digits to the right of the decimal point as there are significant figures in the original number.

2. In an antilogarithm of a number, keep as many digits as there are digits to the right of the decimal point in the original number



• The weak link for *multiplication and division* is the number of *significant figures in the number with the smallest number of significant figures*.



Give the answer of the following operation to the maximum number of significant figures and indicate the number with the greatest relative uncertainty.

$$\frac{35.63 \times 0.5481 \times 0.05300}{1.1689} \times 100\% = 88.5470578\% \\ 88.55\%$$



Calculate the formula weight of  $\text{Ag}_2\text{MoO}_4$  from the individual atomic weights (Ag = 107.870 amu, Mo = 95.94 amu, O = 15.9994 amu); amu = atomic mass unit. Note that the atomic weight of molybdenum is known only to the nearest 0.01 amu, while that for Ag and O are known to 0.001 and 0.0001 amu, respectively. We cannot justifiably say that we know the formula weight of a compound containing molybdenum to any closer than 0.01 atomic unit. Therefore, the most accurately known value for the atomic weight of  $\text{Ag}_2\text{MoO}_4$  is 375.68. All numbers being added or subtracted can be rounded to the least significant unit before adding or subtracting. But for consistency in the answer, one additional figure should be carried out and then the answer rounded to one less figure.

Ag	107.87		0
Ag	107.87		0
Mo	95.94		
O	15.99		94
O	15.99		94
O	15.99		94
O	15.99		94
	<u>375.67</u>		76



It is good practice to keep an extra figure during stepwise calculations and then drop it in the final number.



- **Rounding Off**

If the digit following the last significant figure is greater than 5, the number is rounded, up to the next higher digit. If it is less than 5, the number is rounded to the value of the last significant figure

$$9.47 = 9.5$$

$$9.43 = 9.4$$

If the last digit is a 5, the number is rounded off to the nearest even digit

$$8.65 = 8.6$$

$$8.75 = 8.8$$

$$8.55 = 8.6$$



**There are a number of rules for computations with which the student should be familiar.**

1. Retain as many significant figures in a result or in any data as will give only one uncertain figure. Thus a volume which is known to be between 20.5 mL and 20.7 mL should be written as 20.6 mL, but not as 20.60 mL, since the latter would indicate that the value lies between 20.59 mL and 20.61 mL.
2. In rounding off quantities to the correct number of significant figures, add one to the last figure retained if the following figure (which has been rejected) is 5 or over. Thus the average of 0.2628, 0.2623, and 0.2626 is 0.2626
3. In addition or subtraction, there should be in each number only as many significant figures as there are in the least accurately known number. Thus the Addition  $168.11 + 7.045 + 0.6832$  should be written  $168.11 + 7.05 + 0.68 = 175.84$   
The sum or difference of two or more quantities cannot be more precise than the quantity having the largest uncertainty.
4. In multiplication or division, retain in each factor one more significant figure than is contained in the factor having the largest uncertainty. The percentage precision of a product or quotient cannot be greater than the percentage precision of the least precise factor entering into the calculation. Thus the multiplication  $1.26 \times 1.236 \times 0.6834 \times 24.8652$  should be carried out using the values  $1.26 \times 1.236 \times 0.683 \times 24.87$  and the result expressed to three significant figures.



# Types of errors in experimental data

Determinate Errors  
- **Systematic**  
affect the accuracy of results

Indeterminate Errors—  
**Random**  
(accidental)  
affect measurement precision

## Gross error

An **outlier** is an occasional result in replicate measurements that differs significantly from the other results.

**random (or indeterminate) error**, causes data to be scattered more or less symmetrically around a mean value.

Indeterminate errors are random and cannot be avoided.

**systematic (or determinate) error**, causes the mean of a data set to differ from the accepted value.

Gross errors lead to **outliers**, results that appear to differ markedly from all other data in a set of replicate measurements.



1

## Some common **determinate (systematic)** errors

1. **Instrumental errors.** These include faulty equipment such as uncalibrated glassware. Electronic instruments are also subject to systematic errors.

2. **Operative errors (personal error).** These include personal errors and can be reduced by experience and care of the analyst in the physical manipulations involved. Operative errors can be minimized by having a checklist of operations. Operations in which these errors may occur include transfer of solutions, effervescence and “bumping” during sample dissolution, incomplete drying of samples, and so on. These are difficult to correct for. Other personal errors include mathematical errors in calculations and prejudice in estimating measurements.

3. **Errors of the method. These are the most serious errors of an analysis.** Most of the above errors can be minimized or corrected for, but errors that are inherent in the method cannot be changed unless the conditions of the determination are altered. Some sources of methodical errors include coprecipitation of impurities, slight solubility of a precipitate, side reactions, incomplete reactions, and impurities in reagents. Sometimes correction can be relatively simple, for example, by running a **reagent blank**. A **blank determination** is an analysis on the added reagents only. It is standard practice to run such blanks and to subtract the results from those for the sample.

Determinate or **systematic errors** are nonrandom and occur when something is intrinsically wrong in the measurement.



## Detection of Systematic Instrument and Personal Errors

1. Some systematic instrument errors can be found and corrected by calibration.
2. Most personal errors can be minimized by careful, disciplined laboratory work.
3. Errors due to limitations of the experimenter can usually be avoided by carefully choosing the analytical method or using an automated procedure.

## Detection of Systematic Method Errors

1. analyzing **standard reference materials (SRMs)**, materials that contain **one or more analytes** at known concentration levels.
2. independent and reliable analytical method can be used in parallel with the method being evaluated.
3. *Blank Determinations*. A **blank contains the reagents and solvents used in a determination, but no analyte.**

Determinate errors may be *additive or multiplicative*, depending on the nature of the error or how it enters into the calculation. In order to detect systematic errors in an analysis, it is common practice to add a known amount of standard to a sample (a “spike”) and measure its recovery and note that good spike recovery cannot also correct for response from an unintended analyte (i.e., an interference).

Systematic errors may be either **constant or proportional**.

It is always a good idea to run a blank.

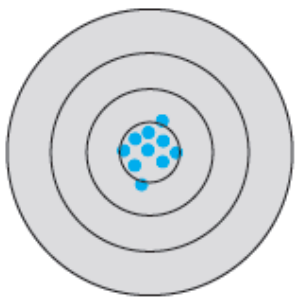


## MINIMISATION OF ERRORS

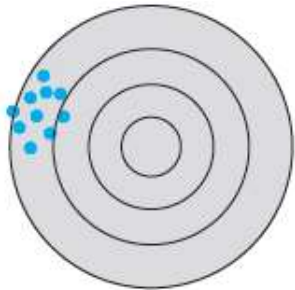
Systematic errors can often be materially reduced by one of the following methods.

- 1. Calibration** of apparatus and application of corrections. All instruments (weights, flasks, burettes, pipettes, etc.) should be calibrated, and the appropriate corrections applied to the original measurements. In some cases where an error cannot be eliminated, it is possible to apply a correction for the effect that it produces; thus an impurity in a weighed precipitate may be determined and its weight deducted.
- 2. Running a blank determination.** This consists in carrying out a separate determination, the sample being omitted, under exactly the same experimental conditions as are employed in the actual analysis of the sample. The object is to find out the effect of the impurities introduced through the reagents and vessels, or to determine the excess of standard solution necessary to establish the end-point under the conditions met with in the titration of the unknown sample. A large blank correction is undesirable, because the exact value then becomes uncertain and the precision of the analysis is reduced.
- 3. Running a control determination.** This consists in carrying out a determination under as nearly as possible identical experimental conditions upon a quantity of a standard substance which contains the same weight of the constituent as is contained in the unknown sample. The weight of the constituent in the unknown can then be calculated. In this connection it must be pointed out that standard samples which have been analysed by a number of skilled analysts are commercially available. These include certain primary standards (sodium oxalate, potassium hydrogenphthalate, arsenic(III) oxide, and benzoic acid) and ores, ceramic materials, irons, steels.
- 4. Use of independent methods of analysis.** In some instances the accuracy of a result may be established by carrying out the analysis in an entirely different manner.
- 5. Running parallel determinations.** These serve as a check on the result of a single determination and indicate only the precision of the analysis.
- 6. Standard addition.** A known amount of the constituent being determined is added to the sample, which is then analysed for the total amount of constituent present. The difference between the analytical results for samples with and without the added constituent gives the recovery of the amount of added constituent. If the recovery is satisfactory our confidence in the accuracy of the procedure is enhanced.

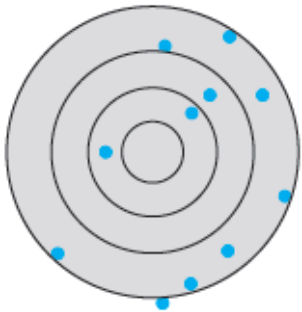




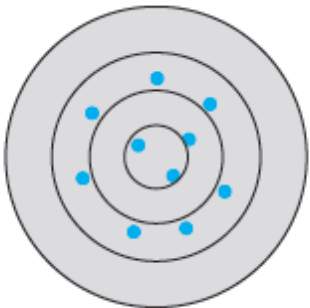
High accuracy, high precision



Low accuracy, high precision



Low accuracy, low precision



High accuracy, low precision

The term error has two slightly different meanings. First, error refers to the difference between a measured value and the "true" or "known" value. Second, error often denotes the estimated uncertainty in a measurement or experiment.

**Precision** describes the reproducibility of measurements—in other words, the closeness of results that have been obtained *in exactly the same way*.

Three terms are widely used to describe the precision of a set of replicate data: **standard deviation**, **variance**, and **coefficient of variation**.

**Accuracy** indicates the closeness of the measurement to the true or accepted value and is expressed by the *error*.

**True Result** - the 'correct' value for a measurement which remains unknown except when a standard sample is being analysed. It can be estimated from the results with varying degrees of precision depending on the experimental method.

Accuracy is expressed in terms of either **absolute or relative error**.

**Replicates** are samples of about the same size that are carried through an analysis in *exactly the same way*.

The precision of a measurement is readily determined by comparing data from carefully replicated experiments. Unfortunately, an estimate of the accuracy is not as easy to obtain. To determine the accuracy, we have to know the true value, which is usually what we are seeking in the analysis.

**Good precision does not guarantee accuracy.**

## Ways of expressing accuracy

the **absolute error** is the difference between an experimental result and an accepted value including its sign.

$$E = x_2 - x_1$$

The **relative error** of a measurement is the absolute error divided by the true value.

$$E_r = \frac{x_2 - x_1}{x_1} \times 100\%$$

The sign of the absolute error tells you whether the value in question is high or low. If the measurement result is low, the sign is negative; if the measurement result is high, the sign is positive.



The results of an analysis are 36.97 g, compared with the accepted value of 37.06 g. What is the relative error?





## Random Errors in Chemical Analysis

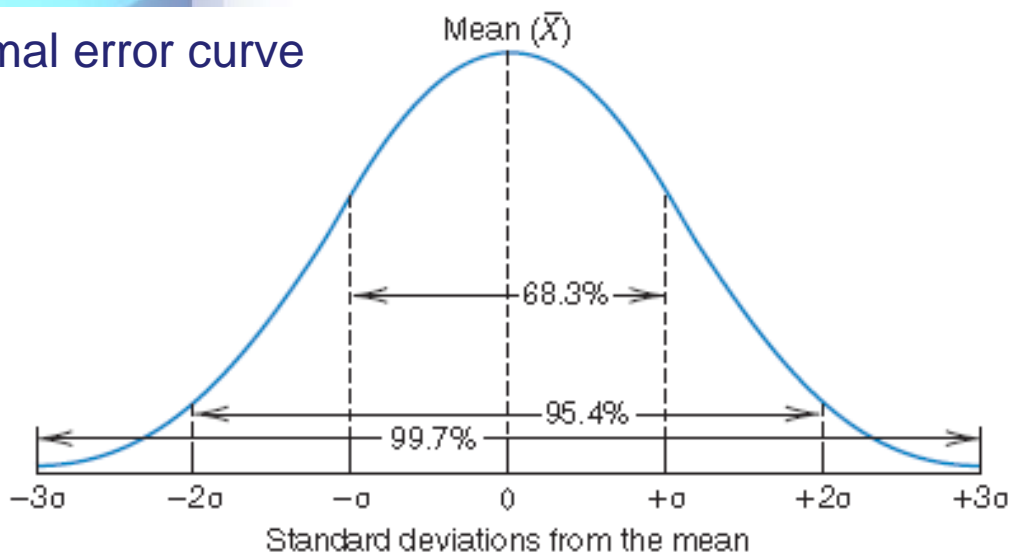
indeterminate errors, often called accidental or **random errors**, which represent the experimental uncertainty that occurs in any measurement. These errors are revealed by small differences in successive measurements made by the same analyst under virtually identical conditions, and they cannot be predicted or estimated. These accidental errors will follow a random distribution.

Indeterminate errors should follow a **normal distribution, or Gaussian curve**.

A **Gaussian, or normal error curve**, is a curve that shows the symmetrical distribution of data around the mean of an infinite set of data

Indeterminate errors really originate in the limited ability of the analyst to control or make corrections for external conditions, or the inability to recognize the appearance of factors that will result in errors.

The symbol  $\sigma$  represents the standard deviation of an infinite population of measurements



## Statistical Treatment of Random Errors

The most widely used measure of central value is the **mean**,  $\bar{x}$ . The mean, also called the **arithmetic mean or the average**, is obtained by dividing the sum of replicate measurements by the number of measurements in the set:

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

$x_i$  - represents the individual values of  $x$  making up the set of  $N$  replicate measurements

The **median** is the middle result when replicate data are arranged in increasing or decreasing order.

$X_{\text{med}}$ , is the middle value when data are ordered from the smallest to the largest value.

To determine the median, we order the data from the smallest to the largest value

3.056 3.080 3.094 3.107 3.112 3.174 3.198

Since there is a total of seven measurements, the median is the fourth value in the ordered data set; thus, the median is 3.107.

**Range** - the numerical difference between the largest and smallest values in a data set ( $w$ ).



**Deviation from the mean  $d_i$**

$$d_i = |x_i - \bar{x}|$$

The **standard deviation  $\sigma$**  of an infinite set of experimental data

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

A statistical measure of the “average” deviation of data from the data’s mean value ( $s$ ).

estimated **standard deviation  $S$**  of a finite set of experimental

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

The **absolute standard deviation,  $S$** , describes the spread of individual measurements about the mean.

**number of degrees of freedom ( $N - 1$ )** The number of independent values on which a result is based

where  $X_i$  is one of  $n$  individual measurements, and  $X$  is the mean. Frequently, the **relative standard deviation,  $S_r$** , is reported.

$$S_r = \frac{s}{\bar{X}}$$

The percent relative standard deviation is obtained by multiplying  $S_r$  by 100%.



## Standard Error of the Mean or Standard deviation of the mean

$$S_m = \frac{s}{\sqrt{N}}$$

The **standard error of a mean**,  $S_m$ , is the standard deviation of a set of data divided by the square root of the number of data points in the set.

The standard deviation is sometimes expressed as the **relative standard deviation (rsd)**, which is just the standard deviation expressed as a fraction of the mean; usually it is given as the *percentage of the mean (% rsd)*, which is often called the **coefficient of variation**.

### Variance ( $s^2$ )

Another common measure of spread is the square of the standard deviation, or the **variance**. The standard deviation, rather than the variance, is usually reported because the units for standard deviation are the same as that for the mean value.

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

The **variance**,  $s^2$ , is equal to the square of the standard deviation.

### The Confidence Limit—How Sure Are You?

$$\text{Confidence limit} = \bar{x} \pm \frac{ts}{\sqrt{N}}$$

$t$  is a statistical factor that depends on the number of degrees of freedom and the confidence level desired.





Calculate the mean and the standard deviation of the following set of analytical results:

15.67, 15.69, and 16.03 g.

$x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
15.67	0.13	0.0169
15.69	0.11	0.0121
16.03	0.23	0.0529
<hr/>	<hr/>	<hr/>
$\Sigma 47.39$	$\Sigma 0.47$	$\Sigma 0.0819$

$$\bar{x} = \frac{\sum x_i}{N} = \frac{47.39}{3} = 15.80$$

$$s = \sqrt{\frac{0.0819}{3 - 1}} = 0.20 \text{ g}$$

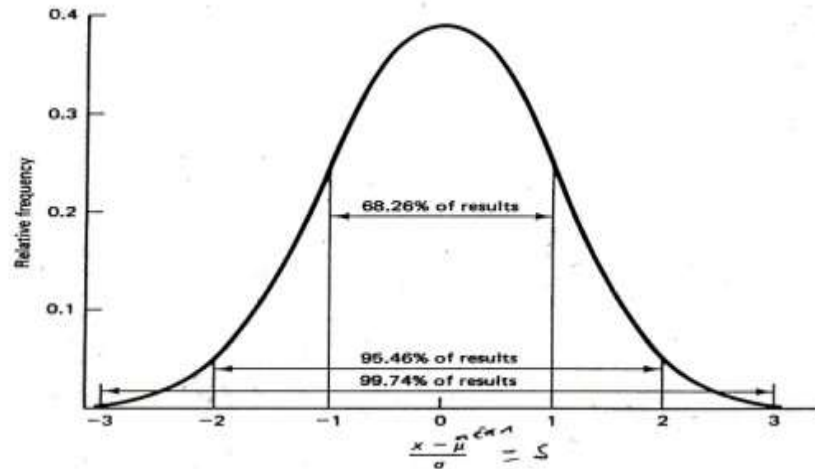
This result would be properly displayed as  $15.8 \pm 0.2 \text{ g}$  (mean  $\pm$  standard deviation) in a final report.



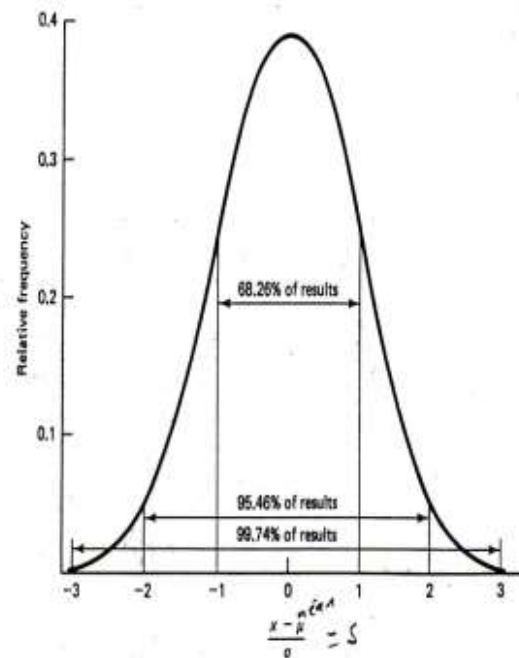
- Standard deviation defines the shape of the normal distribution (particularly width)

- Larger std. dev. more scatter about the mean, worse precision.

- Smaller std. dev. means less scatter about the mean, better precision.

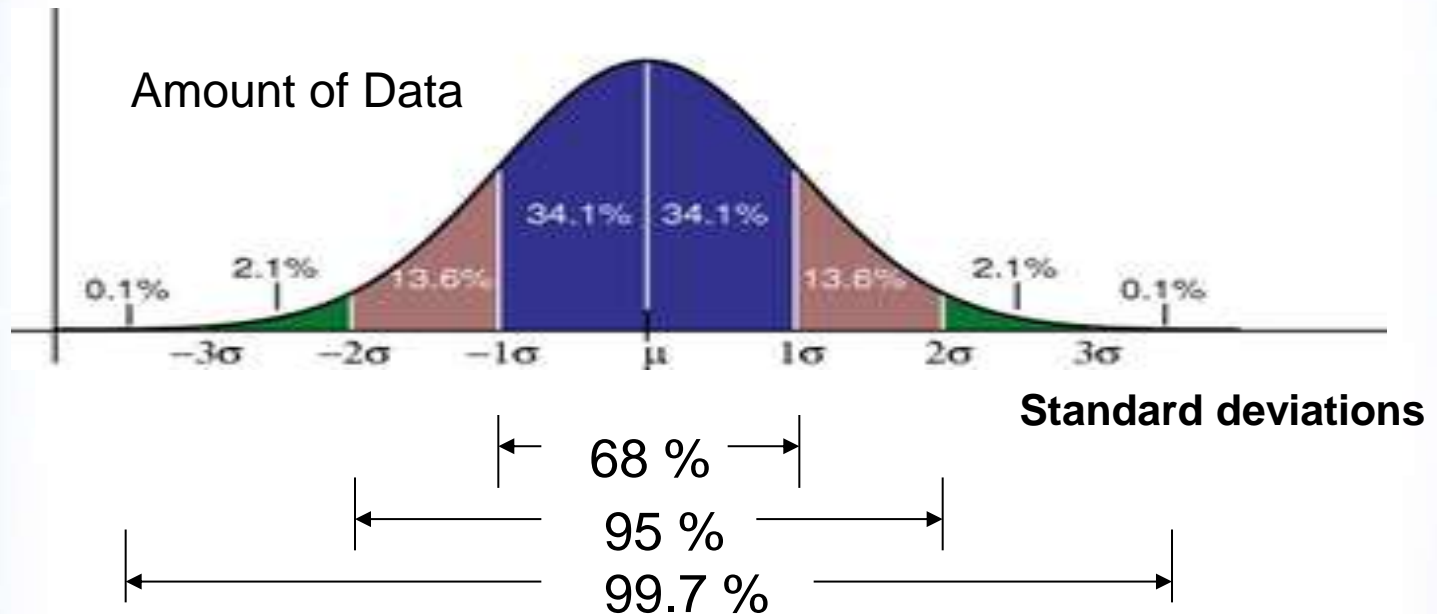


**Figure 2.2** Normal distribution curve; relative frequencies of deviations from the mean for a normally distributed infinite population; deviations  $(x - \mu)$  are in units of  $\sigma$ .



**Figure 2.2** Normal distribution curve; relative frequencies of deviations from the mean for a normally distributed infinite population; deviations  $(x - \mu)$  are in units of  $\sigma$ .

- There is a well-defined relationship between the std. dev. of a population and the normal distribution of the population.
- (May also consider these percentages of area under the curve)





The following replicate weighings were obtained: 29.8, 30.2, 28.6, and 29.7 mg. Calculate the standard deviation of the individual values and the standard deviation of the mean. Express these as absolute (units of the measurement) and relative (% of the measurement) values.

$x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
29.8	0.2	0.04
30.2	0.6	0.36
28.6	1.0	1.00
29.7	0.1	0.01
$\Sigma$ 118.3	$\Sigma$ 1.9	$\Sigma$ 1.41

$$\bar{x} = \frac{118.3}{4} = 29.6$$

$$s = \sqrt{\frac{1.41}{4 - 1}} = 0.69 \text{ mg(absolute)}; \frac{0.69}{29.6} \times 100\% = 2.3\%(\text{coefficient of variation})$$

$$s_{\text{mean}} = \frac{0.69}{\sqrt{4}} = 0.34 \text{ mg(absolute)}; \frac{0.34}{29.6} \times 100\% = 1.1\%(\text{relative})$$







Analyses of a sample of iron ore gave the following percentage values for the iron content: 7.08, 7.21, 7.12, 7.09, 7.16, 7.14, 7.07, 7.14, 7.18, 7.11. Calculate the mean, standard deviation and coefficient of variation for the values.

Results (x)	$x - \bar{x}$	$(x - \bar{x})^2$
7.08	-0.05	0.0025
7.21	0.08	0.0064
7.12	-0.01	0.0001
7.09	-0.04	0.0016
7.16	0.03	0.0009
7.14	0.01	0.0001
7.07	-0.06	0.0036
7.14	0.01	0.0001
7.18	0.05	0.0025
7.11	-0.02	0.0004
$\Sigma x = 71.30$		$\Sigma(x - \bar{x})^2 = 0.0182$
Mean $\bar{x}$ 7.13 per cent		

$$s = \sqrt{\frac{0.0182}{9}}$$

$$= \sqrt{0.0020}$$

$$= \pm 0.045 \text{ per cent}$$

$$\text{C. V.} = \frac{0.045 \times 100}{7.13} = 0.63 \text{ per cent}$$

The mean of several readings will make a more reliable estimate of the true mean than is given by one observation. The greater the number of measurements (n), the closer will the sample average approach the true mean. The standard error of the mean  $S_x$  is given by:

$$s_x = \frac{s}{\sqrt{n}} \quad s_x = \pm \frac{0.045}{\sqrt{10}} = \pm 0.014$$

if 100 measurements were made,

$$s_x = \pm \frac{0.045}{\sqrt{100}} = \pm 0.0045$$

Hence the precision of a measurement may be improved by increasing the number of measurements.



## Statistical Data Treatment and Evaluation

1. Defining a numerical interval around the mean of a set of replicate results within which the population mean can be expected to lie with a certain probability. This interval is called the **confidence interval**. The confidence interval is related to the standard deviation of the mean.
2. Determining the number of replicate measurements required to ensure that an experimental mean falls within a certain range with a given level of probability.
3. Estimating the probability that (a) an experimental mean and a true value or (b) two experimental means are different, that is, whether the difference is real or simply the result of random error. This test is particularly important for discovering systematic errors in a method and determining whether two samples come from the same source.
4. Determining at a given probability level whether the precision of two sets of measurements differs.
5. Comparing the means of more than two samples to determine whether differences in the means are real or the result of random error. This process is known as analysis of variance.
6. Deciding whether to reject or retain a result that appears to be an outlier in a set of replicate measurements.

