

Chapter 8 Fundamental Sampling Distributions and Data Distributions

吳明龍 副教授
成大資訊系/醫資所
minglong.wu@csie.ncku.edu.tw
辦公室:資訊系館 12 樓

8.1 Random Sampling

- Outcome of a statistical experiment
 - Numerical value: total value of a pair of dice tossed
 - **Descriptive representation**: blood types in blood test
- Sampling from distributions or populations
 - Sample mean and sample variance
- The use of high speed computer enhance the use of **formal statistical inference with graphical techniques.**

Random Sampling

- Definition 8.1: A population consists of the totality of the observations with which we are concerned.
 - **Finite size**: 600 students are classified according to blood type
=> a population of size 600
 - **Infinite size**: measuring the atmospheric pressure; some infinite populations are so large
- Each observation in a population is a value of a random variable X having some probability distribution $f(x)$.
 - If one is inspecting items coming off an assembly line for defects, then each observation in population might be a value 0 or 1 of the binomial random variable X with probability distribution
$$b(x;1, p) = p^x q^{1-x}, \quad x = 0,1,$$
where 0 indicates a nondefective item and 1 indicates a defective item.

Random Sampling

- Sometimes, it is impossible or impractical to observe the entire set of observations that make up the population.
- Definition 8.2: A **sample** is a subset of a **population**.
- Inference from the sample to the population are to be valid
 - Obtain **representative samples**
 - **Bias**: Erroneous inferences result from selecting convenient sampling members
 - **Random sample**: independent and at random

Random Sampling

- Definition 8.3: Let X_1, X_2, \dots, X_n be n independent random variables, each having the same probability distribution $f(x)$. We then define X_1, X_2, \dots, X_n to be a random sample of size n from the population $f(x)$ and write its joint probability distribution as

$$f(x_1, x_2, \dots, x_n) = f(x_1)f(x_2)\cdots f(x_n).$$

- If we assume the population of battery lives to be normal, the possible values of any random sample $X_i, i = 1, 2, \dots, 8$, will be precisely the same as those in the original population, and hence X_i has the same identical normal distribution as X .

8.2 Some Important Statistics

- Definition 8.4: Any function of the random variables constituting a random sample is called a statistic.
- Definition 8.5: If X_1, X_2, \dots, X_n represent a random sample of size n , then the sample mean is defined by the statistic

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

- Definition 8.6: If X_1, X_2, \dots, X_n represent a random sample of size n , then the sample variance is defined by the statistic

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

Some Important Statistics

- Example 8.2: A comparison of coffee prices at 4 randomly selected grocery stores in San Diego showed increases from the previous month of 12, 15, 17, and 20 cents for a 1-pound bag. Find the variance of this random sample of price increases.

– Solution

$$\bar{x} = \frac{12 + 15 + 17 + 20}{4} = 16$$

$$\begin{aligned} s^2 &= \frac{\sum_{i=1}^4 (x_i - 16)^2}{4 - 1} = \frac{(12 - 16)^2 + (15 - 16)^2 + (17 - 16)^2 + (20 - 16)^2}{3} \\ &= \frac{34}{3} \end{aligned}$$

Some Important Statistics

- Theorem 8.1: If S^2 is the variance of a random sample of size n , we may write

$$S^2 = \frac{n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2}{n(n-1)}.$$

– Proof

$$\begin{aligned} S^2 &= \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} = \frac{\sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2)}{n-1} \\ &= \frac{\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2}{n-1} = \frac{\sum_{i=1}^n X_i^2 - 2 \left(\frac{\sum_{i=1}^n X_i}{n} \right) \sum_{i=1}^n X_i + n \left(\frac{\sum_{i=1}^n X_i}{n} \right)^2}{n-1} \end{aligned}$$

Some Important Statistics

- (Definition 8.7): The sample standard deviation, denoted by S , is the positive square root of the sample variance.
- Example 8.3: Find the variance of the data 3, 4, 5, 6, 6, and 7, representing the number of trout caught by a random sample of 6 fishermen.

– Solution

$$\sum_{i=1}^6 x_i^2 = 171, \quad \sum_{i=1}^6 x_i = 31$$

$$s^2 = \frac{6 \times 171 - 31^2}{6 \times 5} = \frac{13}{6}$$

8.3 Sampling Distribution

- Statistical inference is concerned with generalizations and predictions.
 - Based on the opinions of several people interviewed on the street, that in a forthcoming election 60% of the eligible voters in the city of Detroit favor a certain candidate.
- Definition 8.5: The probability distribution of a statistic is called a sampling distribution.
 - E.g., the probability distribution of \bar{X} is called the sampling distribution of the mean.
- The sampling distribution of a statistic depends on **the distribution of the population, the size of the samples, and the method of choosing the samples.**

8.4 Sampling Distribution of Means

- Suppose that a random sample of n observations is taken from a normal population with mean μ and variance σ^2 .
- By the reproductive property of the normal distribution established in Theorem 7.11

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

$$\mu_{\bar{X}} = \frac{\mu + \mu + \dots + \mu}{n} = \mu$$

$$\sigma_{\bar{X}}^2 = \frac{\sigma^2 + \sigma^2 + \dots + \sigma^2}{n^2} = \frac{\sigma^2}{n}$$

•Theorem 7.11: If X_1, X_2, \dots, X_n are independent random variables having normal distributions with means $\mu_1, \mu_2, \dots, \mu_n$ and variances $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$, respectively, then the random variable

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$$

has a normal distribution with mean

$$\mu_Y = a_1\mu_1 + a_2\mu_2 + \dots + a_n\mu_n$$

and variance

$$\sigma_Y^2 = a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + \dots + a_n^2\sigma_n^2.$$

Sampling Distribution of Means

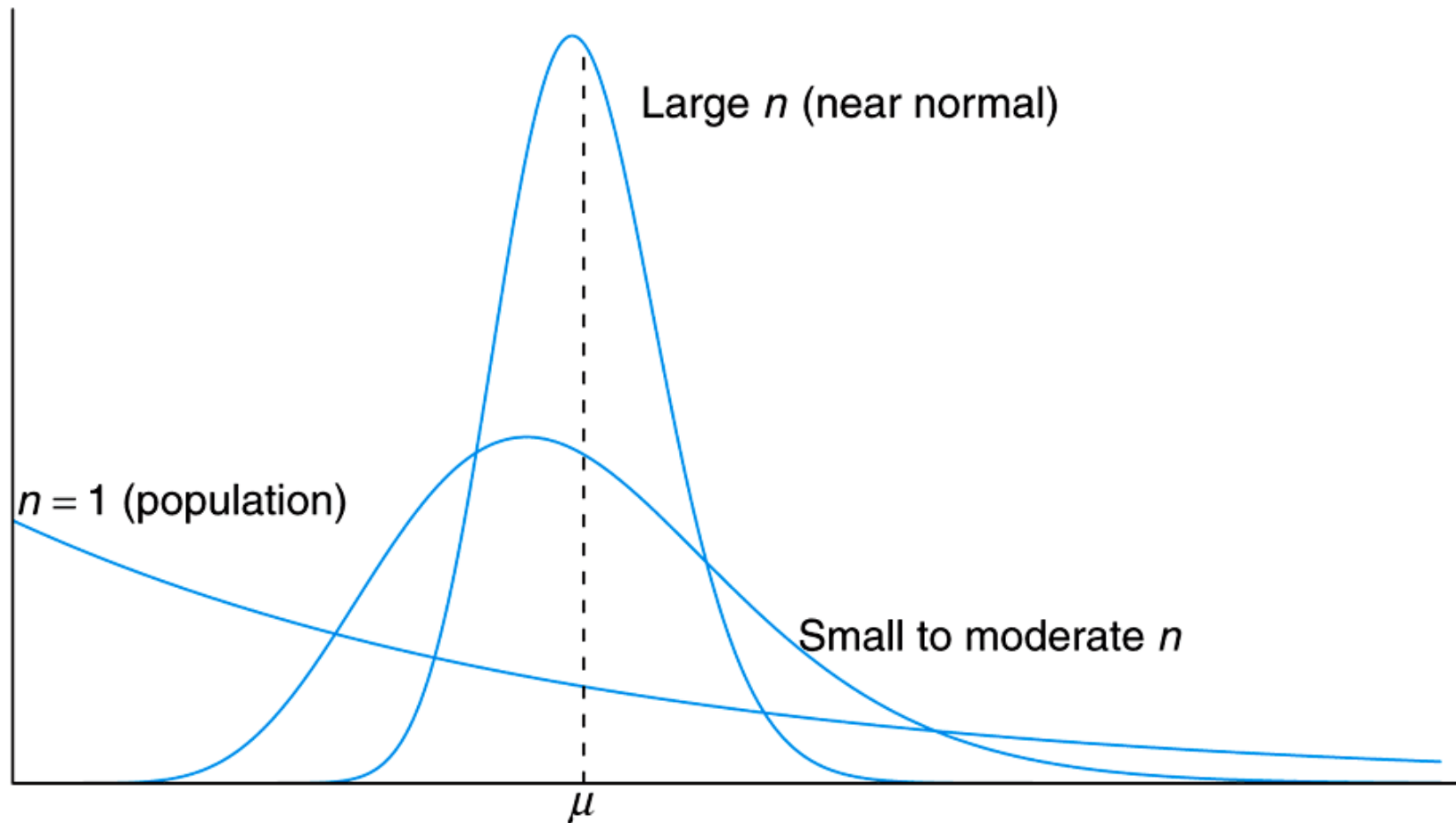
- Theorem 8.2: **(Central Limit Theorem)** If \bar{X} is the mean of a random sample of size n taken from a population with mean μ and finite variance σ^2 , then the limiting form of the distribution of

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}},$$

as $n \rightarrow \infty$, is the standard normal distribution $n(z; 0, 1)$.

- The normal approximation for \bar{X} will generally be good if $n \geq 30$.
- If $n < 30$, the approximation is good only if the population is not too different from a normal distribution.
- If the population is known to be normal, the sampling distribution of \bar{X} will follow a normal distribution exactly, no matter how small the size of the samples.

Figure 8.1 Illustration of the Central Limit Theorem
(distribution of X for $n = 1$, moderate n , and large n)



Sampling Distribution of Means

- Example 8.4: An electric firm manufactures light bulbs that have life mean equal to 800 hours and a standard deviation of 40 hours. Find the probability that a random sample of 16 bulbs will have an average life of less than 775 hours.

– Solution

The sampling distribution of \bar{X} will be approximately normal,

with $\mu_{\bar{X}} = 800, \sigma_{\bar{X}} = 40 / \sqrt{16} = 10$

$$\bar{x} = 775 \Rightarrow z = \frac{775 - 800}{10} = -2.5$$

$$\therefore P(\bar{X} < 775) = P(Z < -2.5) = 0.0062$$

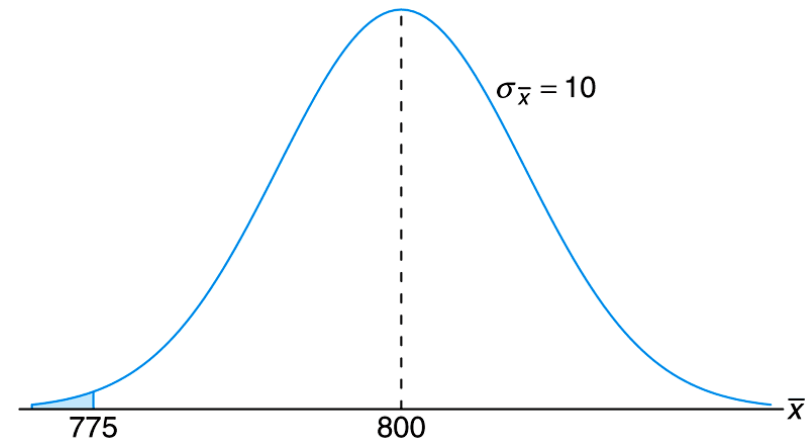


Figure 8.2: Area for Example 8.4

Sampling Distribution of Means

- Case Study 8.1: A engineer conjectures that the population mean of a certain component parts is 5.0 millimeters. An experiment is conducted in which 100 parts produced by the process are selected randomly and the diameter measured on each. It is known that the population standard deviation $\sigma = 0.1$. The experiment indicates a sample average diameter $\bar{x} = 5.027$ millimeters. **Does this sample information appear to support or refute the engineer's conjecture?**

– Solution

$$P[|(\bar{X} - 5)| \geq 0.027] = P[(\bar{X} - 5) \geq 0.027] + P[(\bar{X} - 5) \leq -0.027]$$

$$= 2P\left(\frac{\bar{X} - 5}{0.1/\sqrt{100}} \geq 2.7\right)$$

$$= 2P(Z \geq 2.7) = 2 \times 0.0035 = 0.007.$$

\therefore Strongly refutes the conjecture!

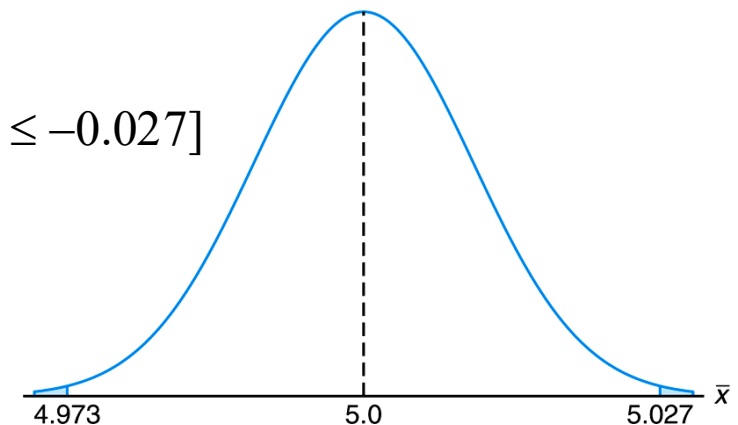


Figure 8.3: Area for Case Study 8.1

Sampling Distribution of the Difference Between Two Averages

- Theorem 8.3: If independent sample of size n_1 and n_2 are drawn at random from two populations, discrete or continuous, with means μ_1 and μ_2 and variances σ_1^2 and σ_2^2 , respectively, then the sampling distribution of the differences of means, $\bar{X}_1 - \bar{X}_2$, is approximately normally distributed with mean and variance given by

$$\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2 \quad \text{and} \quad \sigma_{\bar{X}_1 - \bar{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}.$$

Hence

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{(\sigma_1^2 / n_1) + (\sigma_2^2 / n_2)}}$$

is approximately a standard normal variable.

• Theorem 7.11: If X_1, X_2, \dots, X_n are independent random variables having normal distributions with means $\mu_1, \mu_2, \dots, \mu_n$ and variances $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$, respectively, then the random variable

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$$

has a normal distribution with mean

$$\mu_Y = a_1\mu_1 + a_2\mu_2 + \dots + a_n\mu_n$$

and variance

$$\sigma_Y^2 = a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + \dots + a_n^2\sigma_n^2.$$

Sampling Distribution of the Difference Between Two Averages

- Case Study 8.2: Two independent experiments are being run in which two different types of paints are compared. Eighteen specimens are painted using type *A* and the drying time in hours is recorded on each. The same is done with type *B*. The population standard deviations are both known to be 1.0. Assuming that the mean drying time is equal for the two types of paint, find $P(\bar{X}_A - \bar{X}_B > 1.0)$, where \bar{X}_A and \bar{X}_B are average drying times for samples of size $n_A = n_B = 18$.

– Solution

$$\mu_{\bar{X}_A - \bar{X}_B} = \mu_A - \mu_B = 0 \text{ and } \sigma_{\bar{X}_A - \bar{X}_B}^2 = \frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B} = \frac{1}{18} + \frac{1}{18} = \frac{1}{9}.$$

$$Z = \frac{(\bar{X}_A - \bar{X}_B) - (\mu_A - \mu_B)}{\sqrt{(\sigma_A^2 / n_A) + (\sigma_B^2 / n_B)}} = \frac{1 - 0}{\sqrt{1/9}} = 3.0$$

$$P(Z > 3.0) = 1 - P(Z < 3.0) = 1 - 0.9987 = 0.0013$$

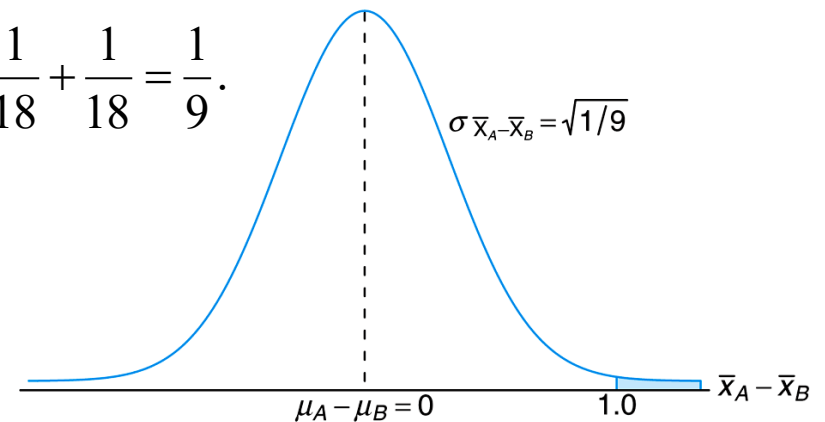


Figure 8.5: Area for Case Study 8.2

Sampling Distribution of the Difference Between Two Averages

- Example 8.6: The television picture tubes of manufacturer *A* have a mean lifetime of 6.5 years and a standard deviation of 0.9 year, while those of manufacturer *B* have a mean lifetime of 6.0 years and a standard deviation of 0.8 year. What is the probability that a random sample of 36 tubes from manufacturer *A* will have a mean lifetime that is at least 1 year more than the mean lifetime of a sample of 49 tubes from manufacturer *B*?

– Solution

$$\mu_{\bar{X}_1 - \bar{X}_2} = 6.5 - 6.0 = 0.5 \text{ and } \sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{0.9^2}{36} + \frac{0.8^2}{49}} = 0.189.$$

$$z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{(\sigma_1^2 / n_1) + (\sigma_2^2 / n_2)}} = \frac{1 - 0.5}{0.189} = 2.65$$

$$\begin{aligned} P(\bar{X}_1 - \bar{X}_2 \geq 1.0) &= P(Z > 2.65) \\ &= 1 - P(Z < 2.65) = 1 - 0.9960 = 0.004 \end{aligned}$$

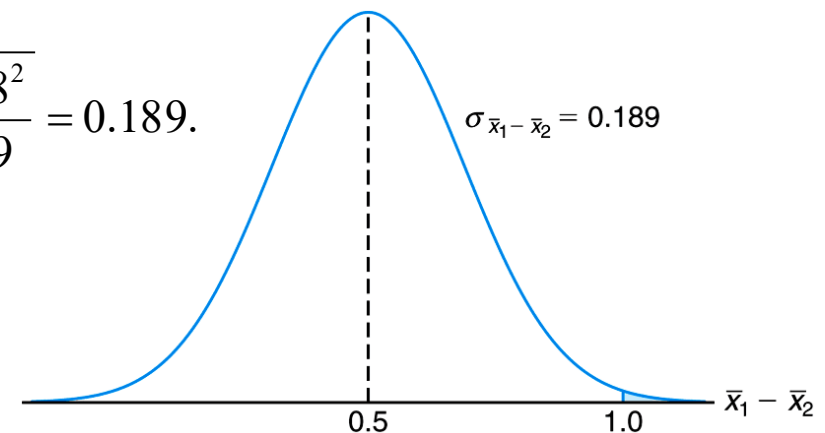


Figure 8.6: Area for Example 8.6

(8.3) Data Displays and Graphical Methods

- Motivation: Use creative displays to extract information about properties of a set.
 - The stem and leaf plots provide the viewer a look at symmetry of the data.
 - Normal probability plots and quantile plots are used to check normal distribution.
- Characterize statistical analysis as the process of drawing conclusion about system variability.
- Statistics provide single measures, whereas a graphical display adds additional information in terms of a picture.

Box and Whisker Plot or Boxplot

- **Box and whisker plot** encloses the interquartile range of the data in a box that has median displayed within.
- **Interquartile range**: between the 75th percentile (upper quartile) and the 25th percentile (lower quartile).
- Boxplot provides the viewer information about outliers which represent “**rare event**”.
- Example 8.3: Nicotine content was measured in a random sample of 40 cigarettes. The data is displayed right.
 - Mild outliers: 0.72, 0.85, and 2.55

1.09	1.92	2.31	1.79	2.28
1.74	1.47	1.97	0.85	1.24
1.58	2.03	1.70	2.17	2.55
2.11	1.86	1.90	1.68	1.51
1.64	0.72	1.69	1.85	1.82
1.79	2.46	1.88	2.08	1.67
1.37	1.93	1.40	1.64	2.09
1.75	1.63	2.37	1.75	1.69

Box and Whisker Plot or Boxplot

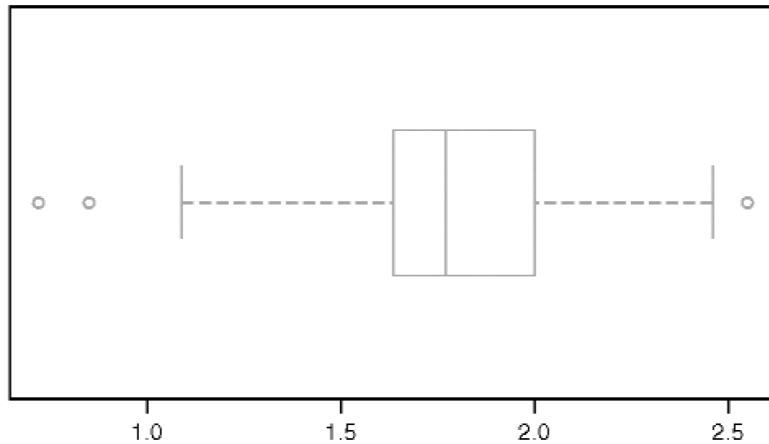


Figure 8.1: Box-and-whisker plot for nicotine data of Exercise 1.21.

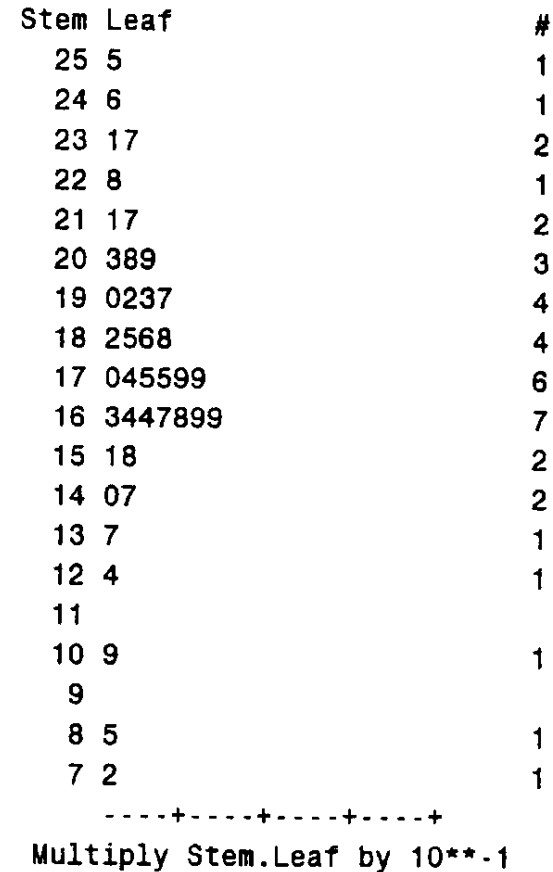


Figure 8.2: Stem-and-leaf plot for nicotine data.

Box and Whisker Plot or Boxplot

- Example 8.4: Consider the following data, consisting of 30 samples measuring the thickness of paint can ears. Figure 8.3 depicts a box and whisker plot for this asymmetric set of data.

Sample	Measurements					Sample	Measurements				
1	29	36	39	34	34	16	35	30	35	29	37
2	29	29	28	32	31	17	40	31	38	35	31
3	34	34	39	38	37	18	35	36	30	33	32
4	35	37	33	38	41	19	35	34	35	30	36
5	30	29	31	38	29	20	35	35	31	38	36
6	34	31	37	39	36	21	32	36	36	32	36
7	30	35	33	40	36	22	36	37	32	34	34
8	28	28	31	34	30	23	29	34	33	37	35
9	32	36	38	38	35	24	36	36	35	37	37
10	35	30	37	35	31	25	36	30	35	33	31
11	35	30	35	38	35	26	35	30	29	38	35
12	38	34	35	35	31	27	35	36	30	34	36
13	34	35	33	30	34	28	35	30	36	29	35
14	40	35	34	33	35	29	38	36	35	31	31
15	34	35	38	35	30	30	30	34	40	28	30

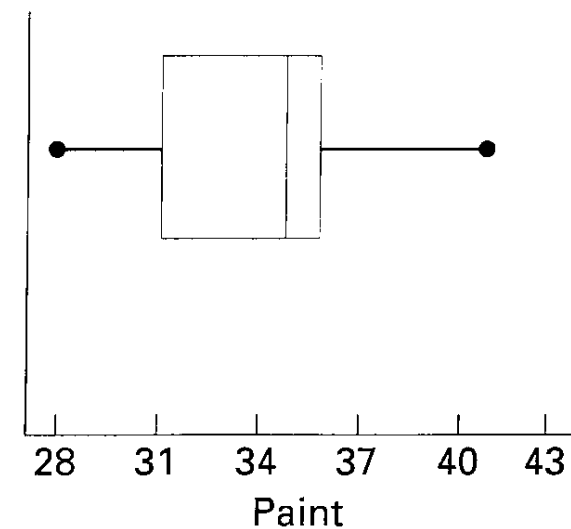


Figure 8.3 Box and whisker plot for thickness of paint can “ears.”

8.8 Quantile and Probability Plots

- Quantile plot
 - Compare samples of data
 - Draw distinctions
 - Depict cumulative distribution function
- Definition 8.6: A quantile of a sample, $q(f)$, is a value for which a specified fraction f of the data values is less than or equal to $q(f)$.
 - Sample median: $q(0.5)$; 75th percentile: $q(0.75)$; 25th percentile: $q(0.25)$;
 - $f_i = \frac{i - \frac{3}{8}}{n + \frac{1}{4}}$, where i is the order of the observations when they are ranked from low to high.

If we denote the ranked observations as $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n-1)} \leq y_{(n)}$
then the quantile plot depicts a plot of $y_{(i)}$ against f_i

Quantile Plot

- In Figure 8.15, quantile plot shows all observations.
 - Large clusters: slopes near zero
 - Sparse data: steeper slopes
- E.g.
 - Sparse data: 28-30
 - High density: 36-38

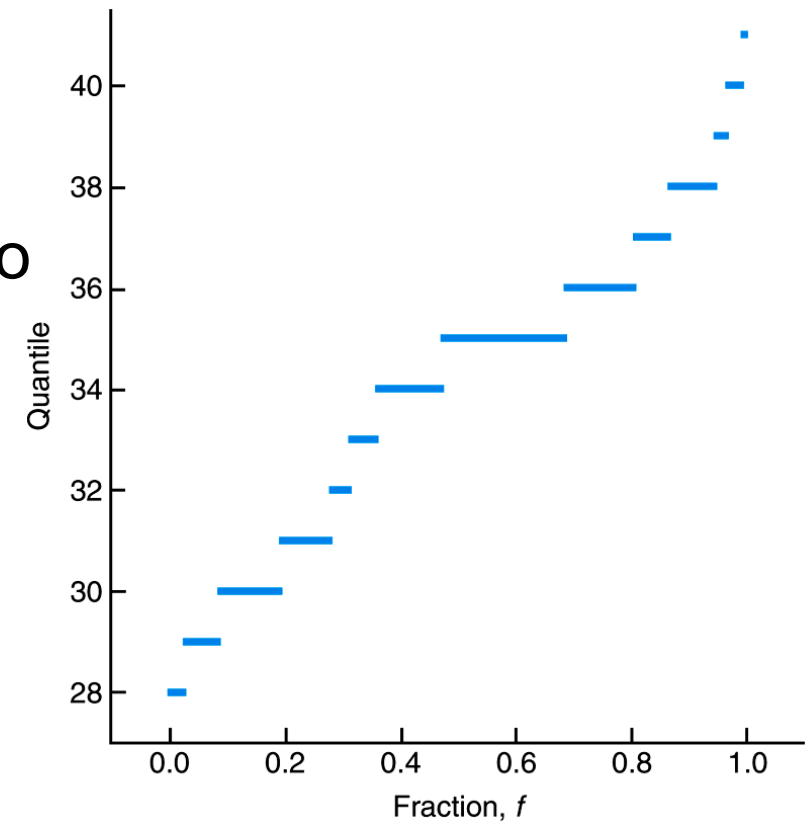


Figure 8.15 Quantile plot for paint data

Normal Quantile-Quantile Plot

- Approximation of quantile of normal distribution

$$q_{\mu,\sigma}(f) = \mu + \sigma \{4.91[f^{0.14} - (1-f)^{0.14}]\}$$

$$q_{0,1}(f) = 4.91[f^{0.14} - (1-f)^{0.14}]$$

- Definition 8.7: The normal quantile-quantile plot is a plot of $y_{(i)}$ (ordered observations)

against $q_{0,1}(f_i)$, where $f_i = \frac{i - \frac{3}{8}}{n + \frac{1}{4}}$.

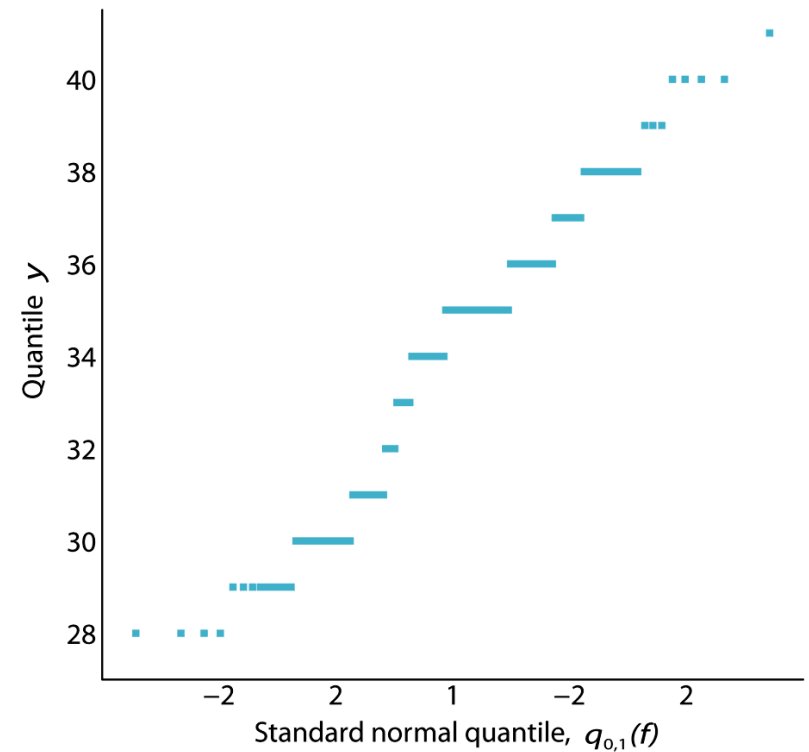


Figure 8.16 Normal quantile-quantile plot for paint data

Normal Quantile-Quantile Plot

- Example 8.12: Construct a normal quantile-quantile plot and draw conclusions regarding whether or not it is reasonable to assume that the two samples are from the same $N(\mu, \sigma)$ distribution.

- Solution

- Far from a straight line
- Station 1 reflect a few values in the lower tail of the distribution and several in the upper tail
- Unlikely

Number of Organisms per Square Meter			
Station 1		Station 2	
5,030	4,980	2,800	2,810
13,700	11,910	4,670	1,330
10,730	8,130	6,890	3,320
11,400	26,850	7,720	1,230
860	17,660	7,030	2,130
2,200	22,800	7,330	2,190
4,250	1,130		
15,040	1,690		

Table 8.1: Data for Example 8.12

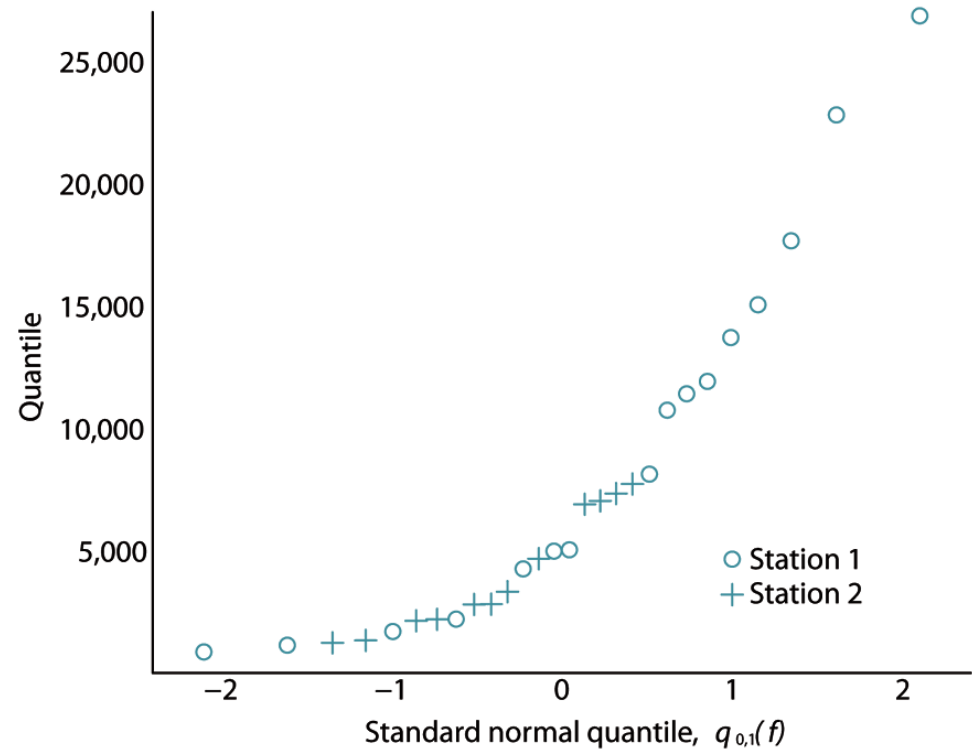


Figure 8.17: Normal quantile-quantile plot for density data of Example 8.12.