

section 5.1

The source of uncertainty in a probability model -

- an uncertainty associated with an outcome or response as described by a probability model. The probability model itself (the probability measure) is fixed. An additional uncertainty comes from the fact that we don't exactly know what this probability measure is. E.g.  $X \sim N(0, \sigma^2)$  - probabilistic uncertainty; in practice we may know (roughly) that  $X$  is Gaussian but  $\sigma^2$  is unknown - thus, only a general model class is known. This is a statistical model.

Statistical inference must be based on data - observed outcomes / responses from a statistical model.

Ex. Stanford Heart Transplant Study (March 1974, JASA)

Can only hope to determine whether a patient lived longer if receiving a new heart transplant - causality is very difficult to establish. Comparisons are made between e.g. the pdf of life length with transplant  $f_T$  and pdf of life length without the transplant  $f_C$ ; T, C - Treatment and Control. The time is measured from the date of determination that the patient is a heart transplant candidate until the end of study. Typically, some characteristics of  $f_T$  and  $f_C$  will be measured e.g. their expectations. Note that clinical vs. statistical significance is still important here.

In practice, only a small number of obs. from  $f_T$  &  $f_C$  are available e.g. 30 patients in a control group and 52 patients in the treatment group. In table 5.1 (control) we know  $X$  and  $S$  - the indicator of whether the patient died by the end of the study.

In the table 5.2, there are additional covariates:

Y - # of days they waited for the transplant

Z - # of days they were alive after the date they received the transplant until the end of the study. Note that  $X = Y + Z$ . Typically, also age, sex etc. may be available as covariates.

## Section 5.2 (Inference Using a Probability Model).

Assume that the probability measure  $P$  is known on a collection of subsets of a sample space  $S$  for a response

$s$ . Typically, we would know  $P$  but uncertain about the future responses  $s \in S$ . Inference about  $s$  may be of interest -

- prediction/estimate of  $s$ ; e.g.  $E[s]$  may be taken as a prediction. Alternatively, we may be asked to construct a subset that has a high probability of containing  $s$  and is in some sense small.

Ex. Lifelength of a machine in years  $X \sim \text{Exp}(1)$  [assumed]

$\Rightarrow E[X] = 1$  - a prediction. The smallest interval containing 95% of all possible values of  $X$  is  $(0, c)$  where

$$0.95 = \int_0^c e^{-x} dx = 1 - e^{-c} \Rightarrow c = -\log(0.05) \approx 2.2957$$

Is  $X_0 = 5$  a plausible value for some machine? Well,

$$P(X > 5) = \int_5^\infty e^{-x} dx = e^{-5} \approx 0.0067.$$

Sometimes, additional prior information is available e.g. that  $s \in C(S)$ .  $\Rightarrow$  Condition  $P(\cdot | C)$  when deriving our inferences. This is a basic axiom  $\rightarrow$  it cannot be proven and may be a source of disagreement sometimes. Principle of conditional probability.

Suppose that  $X \sim \text{Exp}(1)$  and a specific machine has been running for a year. — condition on  $X > 1$ .

The density of the conditional distribution is

$$e^{-(x-1)} \text{ for } X > 1. \quad \text{The predicted lifelength is then}$$

$$\mathbb{E}(X/X > 1) = \int_1^\infty x e^{-(x-1)} dx = 2 \quad \text{— memoryless property}$$

of exponential distribution. The tail probability for  $X_0 = 5$  is  $P(X > 5/X > 1) = \int_5^\infty e^{-(x-1)} dx = e^{-4} =$

$= 0.0183$  — now  $X_0 = 5$  a little more plausible than before.

$$\text{Finally, } 0.95 = \int_1^c e^{-(x-1)} dx = -e^{-(x-1)} \Big|_1^c =$$

$$= 1 - e^{-(c-1)} \Rightarrow e^{-(c-1)} = 0.05 \Rightarrow -(c-1) = \log 0.05, \text{ or}$$
 ~~$c-1 = -\log 20$~~ , or  $c = 1 + \log 20$ .

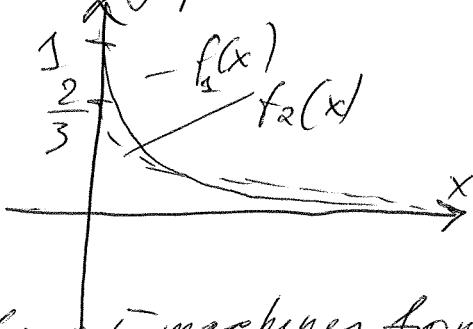
### 5.3 Statistical models.

How to conduct statistical inference for a statistical model  $\{P_\theta : \theta \in \Sigma\}$ ?!  $\theta$  is the parameter of the model that has to be estimated,  $\Sigma$  is the parameter space. Hopefully, the model is identifiable, as  $P_{\theta_1} = P_{\theta_2} \Leftrightarrow \theta_1 = \theta_2$ . If there is a density function, we write  $\{f_\theta : \theta \in \Sigma\}$  instead. Problems we face: 1) to estimate the true parameter value  $\theta$ , 2) construct regions in  $\Sigma$  that contain the true value 3) assess whether the data are in agreement with a suggested value  $\theta_0$ .

For a sample  $X_1, \dots, X_n$  the statistical model for a sample if  $X_1, \dots, X_n$  are independent is  $f_\theta(x_1) f_\theta(x_2) \cdots f_\theta(x_n) \{$  or a sample from the statistical model  $\{f_\theta : \theta \in \Sigma\}$

Distribution:  $X_1 \rightarrow X_n \sim N(\mu, \sigma^2)$ .  $x_1, \dots, x_n$  - their observed values

Two different groups of machines produced by two different manufacturing plants: lifelength  $X \sim \text{Exp}(\theta)$ ,  $\theta \sim \text{Exp}(1.5)$



here  $\text{Exp}(\theta)$  mean  
 $\frac{1}{\theta} e^{-x/\theta}$  so  
 $\mathbb{E}X = \theta$ !

Purchase 5 machines from the same plant - don't know from which one. Observations -  $(x_1, \dots, x_5)$ . State true model as  $f_\theta$ :  $\theta \in \mathbb{R}$ ,  $\mathcal{S} = \{1, 2\}$ . Longer observed lifelengths favor  $\theta=2$ ; i.g. if  $(x_1, \dots, x_5) = (5.0, 3.5, 3.3, 4.1, 2.8)$  we are more certain that  $\theta=2$  than if  $(x_1, \dots, x_5) = (2.0, 2.5, 3.0, 3.1, 3.8)$ . Parameter  $\theta$  is a label but may be an important characteristic of a distribution. ~~Reparametrizations sometimes useful~~ (1  $\leftrightarrow$  1 transformation). Two important examples:

Ex. 5.3.3. Bernoulli model:  $\mathcal{S} = [0, 1]$  if  $x_1 \rightarrow x_n \sim \text{Ber}(\theta)$

$$\Rightarrow f_\theta(x_i) = \theta^{x_i} (1-\theta)^{1-x_i}, \text{ for a sample}$$

$$\prod_{i=1}^n f_\theta(x_i) = \theta^{n\bar{x}} (1-\theta)^{n(1-\bar{x})} \quad \cdot \text{ Can parameterize using } \psi = \log \frac{\theta}{1-\theta} !!$$

Ex. 5.3.4. Location-scale model (Normal)

$(x_1, \dots, x_n) \sim N(\mu, \sigma^2)$  e.g. distribution of heights in a population

$$\Rightarrow \prod_{i=1}^n f_{(\mu, \sigma^2)}(x_i) = (\sigma \sqrt{2\pi})^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\} =$$

$$= (\sigma \sqrt{2\pi})^{-n/2} \exp \left\{ -\frac{n}{2\sigma^2} (\bar{x} - \mu)^2 - \frac{n-1}{2\sigma^2} s^2 \right\}. \text{ derive it here?}$$

Alternative parametrization:  $(\mu, \sigma)$

5.4 Data Collection - draw ~~the~~ the sharp distinction between observational studies and experiments

5.4.1. Finite Populations - a finite set  $N$  of objects (population), a real-valued measurement function  $X$  defined on  $\{1\}$ ;  $X: \{1\} \rightarrow \mathbb{R}$ . E.g.  $X(\pi)$  - height of a student  $\pi$   
 Height - quantitative variable - or  $X(\pi) = \begin{cases} 1 & \text{female} \\ 2 & \text{male} \end{cases}$  - categorical variable. Population and the measurement produce a population distribution ... defined by the population CDF  $F_X: \mathbb{R}^1 \rightarrow [0, 1]$ ,  $F_X(x) = \frac{|\{\pi: X(\pi) \leq x\}|}{N}$ ,

$N = |\{1\}|$ . Consider example 5.4.1. -  $\{1\}$  a population of  $N = 20$  plots of land for the same size;  $X(\pi)$  - a measure of fertility of plot  $\pi$  on a 10-point scale.

Measurements: 4867837546  
 9575834783

$$\Rightarrow F_X(x) = \begin{cases} 0, & x < 3 \\ 3/20, & 3 \leq x < 4 \\ 6/20, & 4 \leq x < 5 \\ 9/20, & 5 \leq x < 6 \\ 13/20, & 6 \leq x < 7 \\ 15/20, & 7 \leq x < 8 \\ 19/20, & 8 \leq x < 9 \\ 1, & 9 \leq x \end{cases}$$

To know  $F_X(x)$  exactly, one must conduct a census - hard!

In practice, we select a subset  $\{\pi_1, \dots, \pi_n\} \subset \{1\}$  where  $n < N$ .

Approximate (estimate)  $F_X(x)$  by the empirical distribution function:  $\hat{F}_X(x) = \frac{|\{\pi_i: X(\pi_i) \leq x, i=1 \dots n\}|}{n}$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{(-\infty, x]}(X(\pi_i)).$$

Note: sometimes multivariate measurements  $X: \{1\} \rightarrow \mathbb{R}^K$  are needed

e.g.  $X(\pi) = (X_1(\pi), X_2(\pi))$

height weight

How should we select the subset  $(\pi_1, \dots, \pi_n)$  and how large  $n$  should be?!

### 5.4.2. Simple random sampling

An important issue is how to select  $\{x_{i_1}, \dots, x_{i_n}\}$  from an entire population in such a way that it does not turn out to have been sampled from a distinct subpopulation? 3) e.g. sampling river creatures from only the top most layer of water. 2) Selecting only students with low ID numbers in order to estimate the population  $F_x$  may result in sampling senior students only. - Called a selection bias (or a selection effect).

This is the main problem with observational studies - the way the data has been generated is unknown to statisticians. Practical solution - select the set  $\{x_{i_1}, \dots, x_{i_n}\}$  ~~selectively~~ randomly. SRS (simple random sampling) - each subset of size  $n$  has the same probability  $\frac{1}{\binom{N}{n}}$  of being selected. In this case,  $(X(\bar{x}_1), \dots, X(\bar{x}_n))$  is random. Thus, for  $n=1$ , we have  $P(X(\bar{x}_1) \leq x) = F_x(x)$ . Go back to the example of  $N=20$  agricultural plots, each one has the prob. of  $\frac{1}{20}$  to be selected. Thus,  $P(X(\bar{x}_1) \leq x) = \frac{P(X(\bar{x}) \leq x)}{20} = F_x(x)$  for every  $x \in \mathbb{R}^1$ .

Prior to observing the sample, we also have  $P(X(\bar{x}_2) \leq x) = F_x(x)$ . However, given that  $X(\bar{x}_1) = x_1$ , we removed one population member with measurement value  $x_1$  so then  $NF_x(x) - 1$  is the # of individuals left in 17 with  $X(\bar{x}) \leq x_1$ . Therefore,

$$P(X(\bar{x}_2) \leq x | X(\bar{x}_1) = x_1) = \begin{cases} \frac{NF_x(x) - 1}{N-1}, & x \geq x_1 \\ \frac{NF_x(x)}{N-1}, & x < x_1 \end{cases}$$

neither result is equal to  $F_x(x)$ .  $\Rightarrow X(\bar{x}_1)$  and  $X(\bar{x}_2)$  are not independent in ~~not~~ a random sample. For large  $N$ , however,  $P(X(\bar{x}_2) \leq x | X(\bar{x}_1) = x_1) \approx F_x(x)$  so  $X(\bar{x}_1)$  and  $X(\bar{x}_2)$  are independent and identically distributed.

### Section 5.4.2.

Similarly, when  $N$  is large and  $n$  is small relative to  $N$ ,  $X(\bar{\alpha}_1), \dots, X(\bar{\alpha}_n)$  are approximately iid with the cdf  $F_X(x)$ . Thus, we will always treat values  $(x_1, \dots, x_n)$  of  $(X(\bar{\alpha}_1), \dots, X(\bar{\alpha}_n))$  as a sample. By the WLLN,  $\hat{F}_X(x) = \frac{1}{n} \sum_{i=1}^n I_{\{X(\bar{\alpha}_i) \leq x\}} \xrightarrow{P} F_X(x)$  as  $n \rightarrow \infty$ .

When the data were collected using SRS, the investigation is called a sampling study. Whenever feasible, they are always preferred to observational studies. Nevertheless, observational studies often have to be used for lack of better opportunities; selection bias always has to be acknowledged. The highest level of statistical evidence is an experiment which is appropriate when cause-effect relationships between variables have to be investigated.

Question #2: how to select the sample size  $n$ ? As large as possible, seemingly... but beware of 1) High costs 2) Increasing possibility of corruption by various errors arising during the collection process. Thus, the best answer - as much as needed for the given accuracy but no more. Therefore, the needed accuracy has to be specified first. These are called sample-size calculations.

If we define  $f_X(x) = \frac{| \{ \bar{\alpha} : X(\bar{\alpha}) = x \} |}{N} = \frac{1}{N} \sum_{\bar{\alpha} \in \Omega} I_{\{X(\bar{\alpha}) = x\}}$ , this  $f_X(x)$  is the proportion of population members that satisfy  $X(\bar{\alpha}) = x$ ; it plays the role of probability function because

$$F_X(x) = \sum_{z \leq x} f_X(z). \text{ This } f_X \text{ is the population relative frequency function.}$$

Based on the sample  $\bar{\alpha}_1, \dots, \bar{\alpha}_n$ ,  $f_X(x)$  may be estimated as  $\hat{f}_X(x) = \frac{| \{ \bar{\alpha}_i : X(\bar{\alpha}_i) = x, i=1, \dots, n \} |}{n} = \frac{1}{n} \sum_{i=1}^n I_{\{X(\bar{\alpha}_i) = x\}}$  - the proportion of sample members satisfying  $X(\bar{\alpha}_i) = x$ . It is appropriate to estimate for categorical variables where  $f_X(x)$  is the proportion in the category specified by  $x$ . With quantitative variables, in general,  $f_X(x)$  is not an appropriate estimation target.

### Section 5.4.3. Histograms.

Now suppose that  $X$  is a continuous quantitative variable.

Group values into intervals:  $(h_1, h_2], (h_2, h_3], \dots, (h_{m-1}, h_m]$

where  $h_1 < h_2 < \dots < h_m$ ,  $(h_1, h_m)$  should cover the range of possible values for  $X$ . Then,

$$h_x(x) = \begin{cases} \frac{|\{x : X(x) \in (h_i, h_{i+1}]\}|}{N(h_{i+1} - h_i)}, & x \in (h_i, h_{i+1}] \\ 0 & \text{otherwise.} \end{cases}$$

Then,  $h_x$  is a density histogram function. This implies that

$h_x(x)(h_{i+1} - h_i)$  is the proportion of individuals for whom  $X(x)$  is in  $(h_i, h_{i+1}]$ . Next,  $F_x(h_j) = \int_{h_i}^{h_j} h_x(x) dx$  for each interval endpoint.

and  $F_x(h_j) - F_x(h_i) = \int_{h_i}^{h_j} h_x(x) dx$ , where  $h_i \leq h_j$ . If the intervals  $(h_i, h_{i+1})$

are small, we expect  $F_x(b) - F_x(a) \approx \int_a^b h_x(x) dx$  for any choice of  $a < b$ .

If the lengths  $h_{i+1} - h_i$  are small and  $N$  is very large,

$h_x$  would be approximated by a smooth continuous function  $f_x(x)$  in the sense that  $\int_a^b f_x(x) dx \approx \int_a^b h_x(x) dx$  for any choice of  $a < b$ . Then, we will

also have  $\int_a^b f_x(x) dx \approx F_x(b) - F_x(a)$ ;  $f_x$  - density function for the

population distribution. This is typically how continuous distributions arise in practice!!

### Section 5.4.4. Survey sampling.

Survey sampling (polling) is entirely based on finite population sampling.

A survey is a set of questions asked of a sample  $\{x_1, \dots, x_n\}$  from a population  $\Pi$ .

If there are  $m$  questions  $\Rightarrow m$  measurements so we have a vector

$$(X_1(\pi), X_2(\pi), \dots, X_m(\pi))$$

- A good example is:
- 1) pre-election polling
  - 2) market surveys done by consumer product companies

Typically, the analysis of results is concerned not just with population distribution of the individual  $X_i$ , but also the joint population distributions. E.g. the joint CDF of  $(X_1, X_2)$  is given by

$$F_{(X_1, X_2)}(x_1, x_2) = \frac{1}{N} \sum_{\delta} I\{\delta : X_1(\delta) \leq x_1, X_2(\delta) \leq x_2\}$$

It is also possible to define  $f_{(X_1, X_2)}(x_1, x_2)$  and joint density histograms if  $(X_1, X_2)$  are both continuous quantitative variables.

Ex 4 mayoral candidates chaotic; 1000 voters randomly selected and asked for whom they vote and their age. The 1st question is

$$X_1(\delta) = 1 - \text{will vote} \quad X_1(\delta) = 0, \text{won't vote. Next, preference is } X_2(\delta) = i, i=1, 2, 3, 4.$$

Typically, we would be interested in joint distributions of  $X_1$  and  $X_2$  but also in joint distributions of  $(X_1, X_3)$  and  $(X_2, X_3)$

Serious problem in survey sampling - how to handle the nonresponse error

# Chapter

- 10 -

## Section 5.5. Some Basic Inferences

### Section 5.5.1. Descriptive Statistics

Ex. 5.5.1.  $f_X(x)$  - proportion of population

members equal to  $x$ ;  $F_X(x)$  - proportion of those  
that don't exceed  $x$ . From a sample  $(x_1, x_2, \dots, x_n)$  we  
can get  $\hat{f}_X(x)$  - proportion of sample values equal to  $x$ .

Then,  $\hat{F}_X(x) = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, x]}(x_i)$  - an empirical cdf.

If  $n=10$ : 3, 2, 2.5, 0.4, 3.3, -2.1, 4.0, -0.3, 2.2, 1.5, 0.2

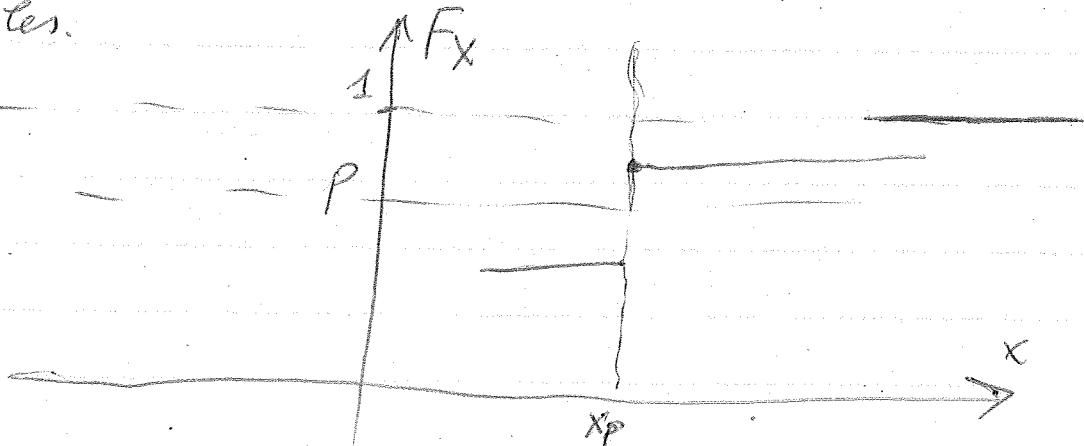
Here,  $\hat{f}_X(x) = 0.1$  for any  $x$  that is a data value and 0  
~~elsewhere~~ elsewhere.  $\hat{F}_X(x) = ?$   $\hat{F}_X(-3) = 0/10 = 0$ ,

$$\hat{F}_X(0) = 2/10 = 0.2, \hat{F}_X(4) = 9/10 = 0.9$$

For  $p \in [0, 1]$ , the  $p$ th quantile (or  $100p$ th percentile)  $x_p$   
is the smallest number  $x_p$  s.t.  $p \leq F_X(x_p)$ . In other  
words,  $x_p = F_X(p) = \min \{x : p \leq F_X(x)\}$ ;

If your test result is at the  $90^{\text{th}}$  percentile,  
your grade is  $x_{0.9}$  and 90% are at or below your  
level. If  $F_X$  is strictly increasing and continuous,

$F_X^{-1}(p)$  is the unique  $x_p$  s.t.  $F_X(x_p) = p$ .  $x_{0.5} = F_X^{-1}(0.5)$  is  
a median,  $x_{0.25} = F_X^{-1}(0.25)$  and  $x_{0.75} = F_X^{-1}(0.75)$  - the 1<sup>st</sup> and  
3<sup>rd</sup> quartiles.



### Section 5.5.1 (continuation).

A natural estimate of a population quantile  $x_p = F_X^{-1}(p)$  is  $\hat{x}_p = F_X^{-1}(p)$  if  $F_X$  is not continuous  $\Rightarrow$  there may not be a

solution to (\*). To solve this issue: 1) Order  $x_1, \dots, x_n$  to obtain the order statistics  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ .

2)  $x_{(i)}$  is the  $(\frac{i}{n})^{th}$  quantile of cdf because

$$F_X(x_{(i)}) = \frac{i}{n} \text{ and } F_X(x) < \frac{i}{n} \text{ when } x < x_{(i)}$$

3) In general, define the sample  $p^{(th)}$  quantile as

$\hat{x}_p = x_{(i)}$  whenever  $\frac{i-1}{n} < p \leq \frac{i}{n}$ . A number of modifications to this are also used, for example, if for some  $(i)$  (\*\*) is satisfied, we can define  $\hat{x}_p = x_{(i-1)} + n(x_{(i)} - x_{(i-1)}) / P - \frac{i-1}{n}$

so that  $\hat{x}_p$  is the linear interpolation between  $x_{(i-1)}$  and  $x_{(i)}$ . When  $n$  is even, this definition gives the sample median as  $\hat{x}_{0.5} = x_{(\frac{n}{2})}$ . An alternative way of defining the sample median is as  $\hat{x}_{0.5} = x_{(\frac{n+1}{2})}$ ,  $n$  is odd.

For  $n$  large enough, the exact choice doesn't matter much.

Here, use the data from Example (5.5.1) using the definition (5.5.3) for  $p=0.5$ ,  $p=0.25$ ,  $p=0.75$ .

### Ex. 5.5.3. Measuring location and scale of a Population Distribution.

Often, we need inferences about the population mean  $\mu_X = \frac{1}{N} \sum X_i$  and the population variance  $\sigma_X^2 = \frac{1}{N} \sum (X_i - \mu_X)^2$ ; here  $N$  is the finite population and  $X_i$  is a real-valued measurement on  $\Omega$ .

For discrete  $X$ , we can also write  $\mu_X = \sum x f_X(x)$  because  $N f_X(x)$  is the

### Section 5.5.1 (continuation)

is the number of elements  $x_i$  with  $X(i) = x$ . In the continuous case,  $P_x \approx \int x f_x(x) dx$  where  $f_x(x)$  is the approximate density.

A natural estimate of the population mean is  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  and the one for  $\sigma_x^2$  is  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ . An alternative

choice would be the population median  $X_{0.5}$  as a measure of location and the population IQR  $X_{0.75} - X_{0.25}$  as a variability measure. This 2<sup>nd</sup> choice is preferred when the distribution is skewed. - an example: house prices in a specific area

The mean & median are very different - example of  $\chi^2$  distribution where the median is 3.3562. In the 2<sup>nd</sup> case, natural estimators are the sample IQR defined as  $\widehat{\text{IQR}} = \widehat{X}_{0.75} - \widehat{X}_{0.25}$ . Again, in the example 5.5.1, check that  $\widehat{X}_{0.5} = 1.5$ ,  $\widehat{\text{IQR}} = 2.75 - 2.05 = 0.70$

### Section 5.5.2 Plotting data.

Using the example 5.5.1, using the same data, show the boxplot using R. Whiskers run from the quartiles to the adjacent values. The adjacent values are given by the greatest value  $\leq$  upper limit ( $3^{\text{rd}} \text{ quartile} + 1.5 \cdot \text{IQR}$ ) - the lower limit ( $3^{\text{rd}} \text{ quartile} - 1.5 \cdot \text{IQR}$ ). Beyond adjacent values are outliers - with  $+$  sign.

Change  $X(0) = 5.0$  to  $X(0) = 15.0$  and you will get an outlier.

Outliers may occur simply because the data has a long-tailed distribution or they may signify a mistake in collecting or recording data.

For categorical variables:	Flavor	Count	Proportion
	Chocolate	42	0.42
	Vanilla	28	0.28
	Butterscotch	22	0.22
	Strawberry	8	0.08

- bars don't need to touch each other.

## Chapter 5.5.3 , Types of Inference.

Descriptive statistics are not that straightforward - we may not know which one we should use. If there's an additional information about a distribution family for the data, e.g. that  $f_X \in \{f_0, f_{0.1}, f_{0.2}\}$ , then descriptive statistics don't use this information at all.

→ using it will help us develop the theory of statistical inference. Recall the 3 types of inference for an unobserved response value  $s$ .

(i) Predict an unknown response value  $s$   
(ii) Construct a subset  $C$  of the sample space  $S$  that has a high probability of containing an unknown response  $s$

(iii) Assess whether or not  $s_0 \in S$  is a plausible value from the probability distribution specified by  $f$ .

In applications, we typically begin with determining some characteristics of the unknown distribution  $f_0$ , e.g. its mean or the median. Denote such a characteristic  $\psi(\theta)$ . E.g. if it is a mean,  $\psi(\theta) = \int x f_0(x) dx$ ; or  $\psi(\theta) = F_0^{-1}(0.5)$ .

→ We consider three types of inference for  $\psi(\theta)$ :

- (i) Choose an estimate  $T(s)$  of  $\psi(\theta)$  - estimation problem
- (ii) Construct a subset  $C(s)$  of the set of possible values for  $\psi(\theta)$  that we believe contains the true value - the problem of credible region (confidence region) construction
- (iii) Assess whether or not  $s_0$  is a plausible value for  $\psi(\theta)$  after having observed  $s$  - hypothesis assessment problem

## Ex. location-scale Normal Model Inference

SRS of heights (in inches) of 30 students:

$$X_1, \dots, X_n \sim N(\mu, \sigma^2), Q = (\mu, \sigma^2) \in \mathcal{Q} = R_x^T R^+$$

Step 1. Use the histogram to check the normality assumption.

Let's say, we are interested in  $\psi(\mu, \sigma^2) = \mu \dots$  or

$$\psi(\mu, \sigma^2) = \bar{x}_{0.90} = \mu + \sigma z_{0.90} \rightarrow 90\text{th percentile of } N(0, 1).$$

By intuition,  $T(x_1, \dots, x_n) = \bar{x}$  is sensible for  $\mu$ :

$$T(x_1, \dots, x_n) = \bar{x} + S z_{0.90} \text{ is sensible for } \mu + \sigma z_{0.90} \text{ later, we'll learn to justify them. If } \bar{x} = 64.577, S = 2.379, z_{0.90} = 1.2816,$$
  
$$\Rightarrow \bar{x} + S z_{0.90} = 64.577 + 2.379(1.2816) = 67.586.$$

To describe the accuracy of  $\bar{x}$  as an est. of  $\mu$ , we will learn to construct a CI of the form  $[\bar{x} - se, \bar{x} + se]$  for some constant  $c$ . . . specifically,  $c = 1.96$  will give a 95% CI for  $\mu$ . The half-length of this interval is  $se = 0.888$  - the measure of the accuracy of  $\bar{x}$  as an est. of  $\mu$ .

Finally, suppose we hypothesized  $\mu_0$  for  $\mu$ . . . says.

$\mu_0 = 65$ . Does it make sense? If  $\bar{x}$  is far from  $\mu_0$ , would seem to be evidence against  $\mu_0$ . . . later, we will base our assessment on

$$z = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} = \frac{64.577 - 65}{2.379/\sqrt{30}} = -1.112 \text{. . . will turn out}$$

to be a plausible value for  $z$ .