

## Chapter 10. Summarizing data

### 1. Empirical probability distribution

IID sample  $(X_1, \dots, X_n)$  with population cdf  $F(x)$

$$\boxed{\text{Empirical cdf } F_n(x) = \text{proportion of } X_i \leq x}$$

For fixed  $x$  sample proportion  $F_n(x)$  is an unbiased and consistent estimate of pop. proportion  $F(x)$

After the sample is collected

$F_n(x)$  is a cdf with mean  $\bar{X}$  and variance  $\frac{n-1}{n}s^2$

### Lifetimes

Lifetime  $T$

cdf  $F(t) = P(T \leq t)$ , pdf  $f(t) = F'(t)$

$$\boxed{\text{Survival function } S(t) = P(T > t) = 1 - F(t)}$$

Empirical survival function  $S_n(t) = 1 - F_n(t)$

the proportion of the data greater than  $t$

$$\boxed{\text{Hazard function } h(t) = f(t)/S(t)}$$

Mortality rate at age  $t$

$$P(t < T \leq t + \delta | T \geq t) \approx \delta \cdot h(t)$$

The negative of the slope of the log survival function

$$h(t) = -\frac{d}{dt} \log S(t)$$

Exp( $\lambda$ ): flat hazard function  $h(t) = \lambda$

$$f(t) = \lambda e^{-\lambda t}, S(t) = e^{-\lambda t}$$

Weibull( $\gamma, \lambda$ ) distribution on  $[0, \infty)$

$$f(t) = \lambda \gamma t^{\gamma-1} e^{-\lambda t^\gamma}, S(t) = e^{-\lambda t^\gamma}, h(t) = \lambda \gamma t^{\gamma-1}$$

scale parameter  $\lambda > 0$  and shape par.  $\gamma > 0$

### Ex 1: Guinea pigs

Guinea pigs infected with tubercle bacillus, p. 349-353

5 treatment and one control group

Fig 10.2: survival function

Fig 10.3: log survival function

### Density estimation

Histogram: observed counts  $O_j$  for cells of width  $h$

small  $h$  - ragged histogram

large  $h$  - obscured histogram, find a balanced  $h$

Scaled histogram

$$f_h(x) = \frac{1}{nh} O_j \text{ for } x \text{ in cell } j \text{ to ensure } \int f_h(x) dx = 1$$

Kernel density estimate with bandwidth  $h$

produces a smooth curve

$$f_h(x) = \frac{1}{nh} \sum \phi\left(\frac{x-X_i}{h}\right), \text{ where } \phi(x) \text{ is the } N(0,1) \text{ pdf}$$

### Ex 2: male heights

If hm is a column of 24 male heights

```
x=160:0.1:210; l=length(x);
```

```
f=normpdf((ones(24,1)*x - hm*ones(1,l))/h);
```

```
fh=sum(f)/(24*h); plot(x,fh)
```

Steam-and-leaf plot for 24 male heights

17:056678899

18:0000112346

19:229

distribution shape plus the numerical information

## 2. Q-Q plots

$p$ -quantile of a distribution  $x_p = F_{-1}(p)$ ,  $0 \leq p \leq 1$

Quantile  $x_p$  cuts off proportion  $p$  of smallest values

$$P(X \leq x_p) = F(x_p) = F(F_{-1}(p)) = p$$

Ordered sample  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$

$$F_n(X_{(k)}) = \frac{k}{n} \text{ and } F_n(X_{(k)} - \epsilon) = \frac{k-1}{n}$$

$X_{(k)}$  is the empirical  $(\frac{k-0.5}{n})$ -quantile

Two samples  $(X_1, \dots, X_n), (Y_1, \dots, Y_m)$

test  $H_0$ : two PDs are equal

by Q-Q plot = plot  $Y$ -quantiles against  $X$ -quantiles

Accept  $H_0$  if the scatter plot is close to the bisector

equal quantiles = equal distributions

Linear model:  $Y = a + b \cdot X$  in distribution

$$P(X \leq x) = P(Y \leq a + bx)$$

Linear model implies linear Q-Q plot  $y_p = a + bx_p$

## Normal probability plot

To test  $H_0: PD = N(\mu, \sigma^2)$  with unspecified parameters

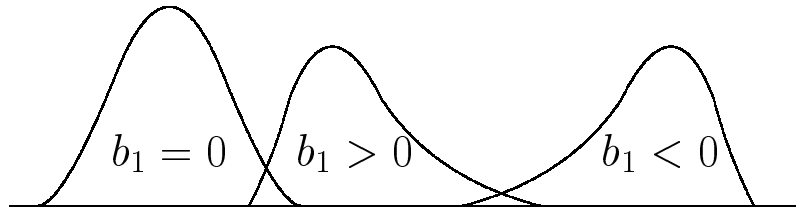
plot the normal quantiles  $\Phi_{-1}(\frac{k-0.5}{n})$  against  $X_{(k)}$

Accept  $H_0$  with  $\mu = a, \sigma = b$

if the scatterplot is close to the line  $x = a + by$

Light tails profile and heavy tails profile

$$\text{Coefficient of skewness: } b_1 = \frac{1}{s^3 n} \sum (X_i - \bar{X})^3$$



$$\text{Kurtosis } b_2 = \frac{1}{s^4 n} \sum (X_i - \bar{X})^4, \text{ normal data } b_2 = 3$$

Leptokurtic distribution:  $b_2 > 3$  heavy tails

platykurtic distribution:  $b_2 < 3$  light tails

## Ex 2: male heights

$$\bar{X} = 181.46, \hat{M} = 180, b_1 = 1.05, b_2 = 4.31$$

Heights of adult males are positively skewed

$$P(\text{height of a random male} < \text{the average}) > 50\%$$

## 3. Measures of location

Central point of a distribution:

population mean  $\mu$ , mode or median  $M$

$M = x_{0.5}$  if distribution is continuous

Population median  $M$ :  $P(X < M) = P(X > M)$

Sample median  $\hat{M} = X_{(k)}$ ,                      if  $n = 2k - 1$   
 $\hat{M} = \frac{X_{(k)} + X_{(k+1)}}{2}$ ,                      if  $n = 2k$

$\hat{M}$  is a robust estimate = insensitive to outliers  
 sample mean  $\bar{X}$  is sensitive to outliers

### Nonparametric sign test

Test  $H_0: M = M_0$  against two-sided  $H_1: M \neq M_0$

sign test statistic:  $Y = \sum I(X_i \leq M_0)$

null distribution  $Y \in \text{Bin}(n, 0.5)$

Reject  $H_0$  if  $M_0$  falls outside  $(X_{(k)}, X_{(n-k+1)})$

where  $k$  is such that  $P(Y < k) = \frac{\alpha}{2}$ ,  $Y \in \text{Bin}(n, 0.5)$

$(X_{(k)}, X_{(n-k+1)}) = \text{nonparametric CI for } M$

$n = 25$ and $k =$	6	7	8	9	10	11	12
$100(1 - \alpha)\%$	99.6	98.6	95.6	89.2	77.0	57.6	31.0

### Trimmed means

Measures of location for the central portion of the data

$\alpha$ -trimmed mean  $\bar{X}_\alpha = \text{sample mean without}$   
 $\frac{n\alpha}{2}$  smallest and  $\frac{n\alpha}{2}$  largest observations

**Ex 2: male heights**  $\bar{X}_{0.4} = 180.36$

When summarizing data compute several  
 measures of location and compare the results

## Nonparametric bootstrap

IID sampling from the empirical distribution

= sampling with replacement from  $x_1, \dots, x_n$

simulate many new samples of size  $n$

Used to view the sampling distribution of an estimate

like trimmed mean, sample median,  $s$

## 4. Measures of dispersion

Sample variance  $s^2$  and sample range  $R = X_{(n)} - X_{(1)}$

are sensitive to outliers

Robust measures of dispersion

interquartile range  $\text{IQR} = x_{0.75} - x_{0.25}$

$\text{MAD} = \text{median of abs dev } |X_i - \hat{M}|, i = 1, \dots, n$

Three estimates of $\sigma$ in $N(\mu, \sigma^2)$ : $s, \frac{\text{IQR}}{1.35}, \frac{\text{MAD}}{0.675}$
------------------------------------------------------------------------------------------------------------

$\text{IQR} = (\mu + \sigma\Phi_{-1}(0.75)) - (\mu + \sigma\Phi_{-1}(0.25)) = 1.35\sigma$

$P\left(\frac{|X-\mu|}{\sigma} \leq z\right) = 0.5$  for  $z = \Phi_{-1}(0.75) = 0.675$

## Boxplot

box center = median

upper edge of the box = upper quartile (UQ)

lower edge of the box = lower quartile (LQ)

upper whisker end =  $\{\text{max data point} \leq \text{UQ} + 1.5 \text{ IQR}\}$

lower whisker end =  $\{\text{min data point} \geq \text{LQ} - 1.5 \text{ IQR}\}$

dots =  $\{\text{data} \geq \text{UQ} + 1.5 \text{ IQR}\}$  and  $\{\text{data} \leq \text{LQ} - 1.5 \text{ IQR}\}$

Convenient to compare different samples

Fig 10.14, p.374: daily  $\text{SO}_2$  concentration data