

Chapter 10. Summarizing data

1 Empirical probability distribution

Population cumulative distribution function $F(x) = P(X \leq x)$. For an IID sample (X_1, \dots, X_n) define

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

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

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

Lifelength T with cdf $F(t) = P(T \leq t)$ and pdf $f(t) = F'(t)$.

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

Empirical survival function $S_n(t) = 1 - F_n(t)$ is the proportion of the data greater than t .

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

Mortality rate at age t : as δ tends to zero, $P(t < T \leq t + \delta | T \geq t) = \frac{F(t+\delta) - F(t)}{S(t)} \sim \delta \cdot h(t)$.

It is also the negative of the slope of the log survival function: $h(t) = -\frac{d}{dt} \log S(t) = -\frac{d}{dt} \log(1 - F(t))$.

Example. 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.

The flat hazard function $h(t) = \lambda$ corresponds to the $\text{Exp}(\lambda)$ model with $f(t) = \lambda e^{-\lambda t}$ and $S(t) = e^{-\lambda t}$. Weibull(γ, λ) distribution with the shape parameter $\gamma > 0$:

$$f(t) = \lambda \gamma t^{\gamma-1} e^{-\lambda t^\gamma}, t \geq 0, S(t) = e^{-\lambda t^\gamma}, h(t) = \lambda \gamma t^{\gamma-1}$$

2 Density estimation

Histogram: plot observed counts O_j for cells of width h . Small h - ragged histogram, large h - obscured histogram, find a balanced h .

Scaled histogram: plot $f_h(x) = \frac{1}{nh} O_j$ for x in cell j to ensure $\int f_h(x) dx = 1$.

Kernel density estimate with bandwidth h produces a smooth curve

$$\boxed{f_h(x) = \frac{1}{nh} \sum \phi\left(\frac{x - X_i}{h}\right), \text{ where } \phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

Example. Male heights.

If hm is a column of 24 male heights, the for a given bandwidth h the following matstat code produces a plot for the kernel density estimate

```
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)
```

Stem-and-leaf plot for 24 male heights indicates the distribution shape plus gives the numerical information:

17:056678899
 18:0000112346
 19:229

3 Q-Q plots

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

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

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

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

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

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

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

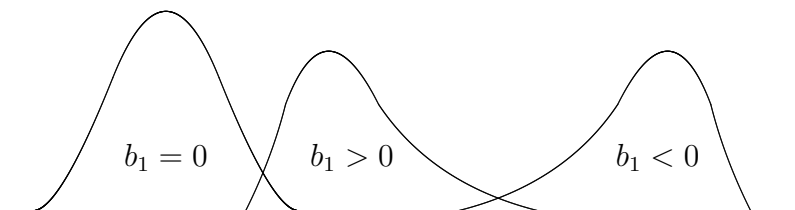
Linear model implies linear Q-Q plot $y_p = a + bx_p$

Normal probability plot. To test visually the normality hypothesis H_0 : PD = $N(\mu, \sigma^2)$ with unspecified parameters plot the normal quantiles $\Phi_{-1}\left(\frac{k-0.5}{n}\right)$ against $X_{(k)}$.

Accept H_0 with $\mu = a, \sigma = b$, if the scatterplot is close to the straight line $x = a + by$.

If normality does not hold, draw a straight line via empirical lower and upper quartiles to detect a light tails profile or 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).

Example. Male heights. Summary statistics: $\bar{X} = 181.46, \hat{M} = 180, b_1 = 1.05, b_2 = 4.31$. Heights of adult males are positively skewed: $P(\text{height of a random male} < \text{the average}) > 50\%$.

4 Measures of location

Central point of a distribution: either population mean μ , or mode, or median M defined as $M = x_{0.5}$, if distribution is continuous.

$$\boxed{\text{Population median } M: P(X < M) = P(X > M)}$$

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

The sample median \hat{M} is a robust estimate, that is insensitive to outliers, while the sample mean \bar{X} is sensitive to outliers.

Nonparametric sign test

Given an iid sample test $H_0: M = M_0$ against the two-sided alternative $H_1: M \neq M_0$. No parametric model is assumed. The sign test statistic $Y = \sum_{i=1}^n I(X_i \leq M_0)$ counts the number of observations below the null hypothesis value. Under the null hypothesis

$$P(X_{(k)} < M_0 < X_{(n-k+1)}) = P(k \leq Y \leq n - k),$$

and $Y \sim \text{Bin}(n, 0.5)$.

$$\boxed{(X_{(k)}, X_{(n-k+1)}) = \text{nonparametric } 100 \cdot P(k \leq Y \leq n - k)\% \text{ CI for the population median.}}$$

Reject H_0 if M_0 falls outside the corresponding confidence interval $(X_{(k)}, X_{(n-k+1)})$.

Example. For $Y \in \text{Bin}(n, 0.5)$ with $n = 25$ we have

for $k =$	6	7	8	9	10	11	12
$P(k \leq Y \leq n - k) =$	99.6	98.6	95.7	89.2	77.0	57.6	31.0

Thus $(X_{(8)}, X_{(18)})$ is a 95.7% CI for the median.

Trimmed means

Measures of location for the central portion of the data

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

Example. Male heights. Ignoring 20% of largest and 20% of smallest observations we compute $\bar{X}_{0.4} = 180.36$.

$$\boxed{\text{When summarizing data compute several measures of location and compare the results}}$$

Nonparametric bootstrap

IID sampling from the empirical distribution = sampling with replacement from x_1, \dots, x_n .

Simulate many new samples of size n to get an idea of the sampling distribution of an estimate like trimmed mean, sample median, s .

5 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 $\text{IQR} = x_{0.75} - x_{0.25}$

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

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

In the $N(\mu, \sigma^2)$ case $\text{IQR} = (\mu + \sigma\Phi_{-1}(0.75)) - (\mu + \sigma\Phi_{-1}(0.25)) = 1.35\sigma$, because $\Phi_{-1}(0.75) = 0.675$.
Moreover, $\text{MAD} = 0.675\sigma$, since $P(|X - \mu| \leq 0.675\sigma) = 0.5$.

Boxplot

box center = median

upper edge of the box = upper quartile (UQ)

lower edge of the box = lower quartile (LQ)

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

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

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. See for example Fig 10.14, p.374: daily SO_2 concentration data.