Statistical inference (MVE155/MSG200)

Summarising data

<ロト < 回 ト < 巨 ト < 巨 ト < 巨 ト 三 の Q () 1/22

- Estimating distribution function
- Quantiles and quantile-quantile plots
- Estimating density function
- Skewness and kurtosis
- Measures of dispersion
- Boxplot

Distribution function F of a random variable X is defined as

$$F(x) = P(X \le x) = \begin{cases} \int_{-\infty}^{x} f(y) \, dy \\ \sum_{y \le x} P(X = y) \end{cases}$$

The empirical distribution function based on the sample $(x_1, ..., x_n)$ is

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(x_i \leq x),$$

where \mathbb{I} is an indicator function, is an unbiased and consistent estimator for F(x).

Empirical distribution function

As a function of x, $\hat{F}(x)$ is a distribution function for a uniform random variable Y with

$$P(Y = x_i) = \frac{1}{n}, \quad i = 1, ..., n$$

if $x_i \neq x_j$ for all $i \neq j$. We have that (even when some of the x_i 's coinside)

$$\mathbb{E}(Y) = \sum_{i=1}^n x_i \operatorname{P}(Y = x_i) = \sum_{i=1}^n \frac{x_i}{n} = \bar{x}$$

and

$$Var(Y) = \mathbb{E}(Y^{2}) - (\mathbb{E}(Y))^{2}$$

= $\sum_{i=1}^{n} x_{i}^{2} P(Y = x_{i}) - \left(\sum_{i=1}^{n} x_{i} P(Y = x_{i})\right)^{2}$
= $\overline{x^{2}} - (\overline{x})^{2} = \frac{n-1}{n}s^{2}$

and $\frac{n-1}{n}s^2$ is called empirical variance.

< □ > < ⑦ > < 言 > < 言 > 言 少 < ⊙ 4/22 When studying life length L, survival function

 $S(x) = P(L > x) = 1 - P(L \le x) = 1 - F(x), \quad x \ge 0$

and its empirical counterpart

$$\hat{S}(x) = 1 - \hat{F}(x)$$

are often used since it gives you the probability of living longer than a certain age x.

Mortality rate at age x can be defined as f(x)/S(x) since

$$\frac{\mathrm{P}(x < L \le x + \delta | L \ge x)}{\delta} \to \frac{f(x)}{S(x)}$$

as $\delta \rightarrow 0$. This is called a hazard rate, i.e.

$$h(x) = \frac{f(x)}{S(x)} = \frac{f(x)}{1 - F(x)}.$$

For the exponential distribution, the hazard rate is constant.

Hazard rate is also the negative slope of the log survival function since

$$-\frac{d}{dx}\ln(S(x)) = -\frac{S'(x)}{S(x)} = \frac{f(x)}{S(x)}$$

Quantiles

The inverse of the distribution function F, F^{-1} , is called a quantile function

$$Q({m p}) = {m F}^{-1}({m p}), \quad 0 < {m p} < 1$$

and the p quantile is defined as

 $x_p = Q(p).$

For an ordered sample $(x_{(1)}, ..., x_{(n)})$, where $x_{(1)} = \min\{x_1, ..., x_n\}$ and $x_{(n)} = \max\{x_1, ..., x_n\}$,

$$\hat{F}(x_{(k)}) = rac{k}{n}$$
 and $\hat{F}(x_{(k)} - \epsilon) = rac{k-1}{n}$

(where $\epsilon > 0$ small). Therefore, $x_{(k)}$ is called the empirical $\frac{k-0.5}{n}$ -quantile.

- median $m = x_{0.5} = Q(0.5)$
- lower quartile $x_{0.25} = Q(0.25)$
- upper quartile $x_{0.75} = Q(0.75)$

We have two independent equal sized samples, $(x_1, ..., x_n)$ with the distribution function F_1 and quantile function Q_1 and $(y_1, ..., y_n)$ with the distribution function F_2 and quantile function Q_2 . We want to test whether these two samples come from the same distribution, i.e. we test

$$H_0:F_1\equiv F_2$$

or equivalently,

 $H_0: Q_1 \equiv Q_2.$

Graphically, this can be done by a QQ plot which is a scatter plot with coordinates $(x_{(k)}, y_{(k)})$, k = 1, ..., n.

If the QQ plot is approximately a 45° line, the distributions are approximatively the same (and H_0 is not rejected).

If the QQ plot is a straight line with some other angle, X and Y have a linear relationship

Y = a + bX

meaning that

 $F_1(x) = F_2(a + bx)$

and

 $Q_2(p) = a + bQ_1(p).$

To test whether the data can be considered being normally distributed, we define

$$y_k = \Phi^{-1}\left(\frac{k-0.5}{n}\right), \quad k = 1, ..., n,$$

where Φ^{-1} is the quantile function for the standard normal distribution N(0, 1). Then, plot

 $(x_{(1)}, y_{(1)}), ..., (x_{(n)}, y_{(n)})$

and check whether the QQ-plot is a straight line.

Normal QQ-plot: heavy tails



Red curve: N(0, 1)Black curve: distribution with heavy tails



Normal QQ-plot: light tails



Red curve: N(0, 1)Black curve: distribution with light tails



・ロ ・ (日 ト く 三 ト く 三 ト く 三) 2 の 2 で
13/22

Median is the 0.5 quantile, $m = x_{0.5} = Q(0.5)$, and

$$P(X < m) = P(X > m) = 0.5.$$

It can be estimated from the ordered sample $x_{(1)}, ..., x_{(n)}$ by

$$\hat{m} = \begin{cases} x_{(k)} & \text{if } n = 2k - 1 \pmod{\frac{1}{2}(x_{(k)} + x_{(k+1)})} & \text{if } n = 2k \pmod{\frac{1}{2}(x_{(k)} + x_{(k+1)})} \end{cases}$$

Confidence interval for median

Let us have a sample $(x_1, ..., x_n)$ (without ties) from a continuous population distribution and let

$$Y=\sum_{i=1}^{n}\mathbb{I}(x_{i}\leq m).$$

Then, $Y \sim Bin(n, 0.5)$ and a $100(1 - 2p_k)\%$ confidence interval for *m* becomes

$$I_m = (x_{(k)}, x_{(n-k+1)}),$$

where

 $p_k = \mathrm{P}(Y < k).$

Y is also used as a test statistic in the sign test, which can be used instead of the one sample t-test when the sample size is small and the data not normal. Reject the $H_0: m = m_0$ in the favour of $H_1: m \neq m_0$ if the value of Y is not within the confidence interval.

Density estimation

A density function can be estimated by using the histograms of the observed counts

$$c_j = \sum_{i=1}^n \mathbb{I}(x_i \in \operatorname{cell}_j),$$

where the observation interval is divided into adjacent cells of width h. The density function can be estimated by the scaled histogram

$$f_h(x) = \frac{c_j}{nh}$$
 for $x \in \operatorname{cell}_j$.

Often, the scaled histogram is smoothened resulting in a kernel estimate with bandwidth h:

$$f_h(x) = rac{1}{nh} \sum_{i=1}^n \phi\left(rac{x-x_i}{h}
ight),$$

where (for example) $\phi(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$.

16 / 22



Kernel estimation: selection of bandwidth



Skewness and kurtosis

Let the random variable X come from a distribution with mean μ and variance σ^2 . For the standardized version of X,

$$Z=\frac{X-\mu}{\sigma},$$

the first moment (mean) $\mu_1 = \mathbb{E}(Z) = 0$ and the second moment $\mu_2 = \mathbb{E}(Z^2) = 1$.

Population coefficient for skewness and population kurtosis are defined as

$$\mu_3=\mathbb{E}(Z^3)$$
 and $\mu_4=\mathbb{E}(Z^4),$

which, given \bar{x} and s, can be estimated by

$$m_3 = rac{1}{ns^3} \sum_{i=1}^n (x_i - ar{x})^3$$
 and $m_4 = rac{1}{ns^4} \sum_{i=1}^n (x_i - ar{x})^4.$

イロト 不得 トイヨト イヨト ヨー ろんの

For normal distribution, $\mu_3 = 0$ (symmetric) and $\mu_4 = 3$ (reference, not heavy or light tails).

Skewness

- $\mu_3 = 0$, symmetric
- $\mu_3 > 0$, skewed to the right
- $\mu_3 < 0$, skewed to the left

Kurtosis

- $\mu_4 > 3$, heavy tails
- $\mu_4 < 3$, light tails

- Sample variance and standard deviation: S² and S
- Sample range: $x_{(n)} x_{(1)}$
- ▶ Interquartile range, IQR: $x_{0.75} x_{0.25}$
- ▶ Median of the absolute values of deviations, MAD: sample median of {|x_i m̂|, i = 1, ..., n}

- Box: from the lower quartile to the upper quartile with the median indicated with the thicker line.
- Whiskers: Based on the 1.5 IQR value.
- Outliers: Observations outside the whiskers.

