

Statistics

Statistics - Basic

- *mode* (眾數): the value with highest frequency
- *median* (中位數): the value in the middle
- *mean* (平均數): the average value ($\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$)

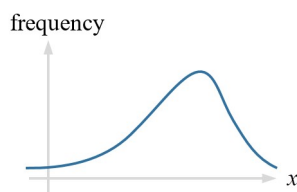
Statistics - Basic

Example. 23, 29, 20, 32, 23, 33, 25, 21

X	20	21	23	25	29	32	33
frequency	1	1	2	1	1	1	1

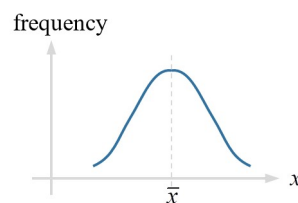
- $\text{mode}(23, 29, 20, 32, 23, 33, 25, 21) = 23$
- $\text{median}(23, 29, 20, 32, 23, 33, 25, 21) = \text{median}(20, 21, 23, 23, 25, 29, 32, 33) = \frac{23+25}{2} = 24$
- $\text{mean}(20, 21, 23, 23, 25, 29, 32, 33) = \frac{23+29+20+32+23+33+25+21}{8} = 25.75$

Skewness



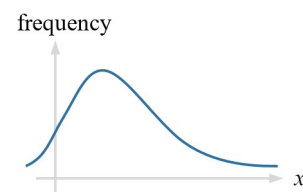
(a) Negative skew:

$$x_{\text{median}} < x_{\text{mode}}$$



(b) Symmetric:

$$\bar{x} = x_{\text{median}} = x_{\text{mode}}$$



(c) Positive skew:

$$x_{\text{mode}} < x_{\text{median}}$$

Statistical Graphs

Statistical Graphs (a.k.a. statistical diagrams, 統計圖表) = {bar chart (長條圖), histogram (直方圖), scatter plot (散佈圖), line chart (折線圖), boxplot (箱形圖), pie chart (圓形圖), ...}

Expectation

X	x_1	x_2	\dots	x_m
$frequency$	f_1	f_2	\dots	f_m

$$p_i = \frac{f_i}{\sum_{i=1}^m f_i} \Rightarrow$$

X	x_1	x_2	\dots	x_m
P	p_1	p_2	\dots	p_m

Expectation (期望值, a.k.a. expected value),

$$E(X) = \frac{\sum_{i=1}^m x_i f_i}{\sum_{i=1}^m f_i} = \sum_{i=1}^m x_i P(X = x_i) = \sum_{i=1}^m x_i p_i \quad \text{where } p_i \equiv p(x_i) \equiv P(X = x_i) \equiv \frac{f_i}{\sum_{i=1}^m f_i}$$

$$\Rightarrow E(X) = \sum_{i=1}^m x_i p_i \quad \text{where } \sum_{i=1}^m p_i = 1$$

Expectation

Properties:

(1) $E(aX + b) = aE(X) + b$

Proof:

$$E(aX + b) = \sum_{i=1}^m (ax_i + b)P(x_i) = a \sum_{i=1}^m x_i p(x_i) + b \sum_{i=1}^m p(x_i) = aE(X) + b$$

□

Bias of Estimates (估計之偏差)

Let $\hat{\theta}$: estimate of θ ($\hat{\theta} \in \Theta$)

where $\hat{\theta}$ is called the *estimate* (估計值), and

Θ is called the *estimator* (估計量),

the *bias of estimate* (偏誤估計) is defined as

$$\text{bias}(\hat{\theta}) \equiv E(\hat{\theta}) - \theta.$$

$\hat{\theta}$ is called

(i) *unbiased*, if $\text{bias}(\hat{\theta}) = 0$;

(ii) *biased*, if $\text{bias}(\hat{\theta}) \neq 0$;

(iii) *asymptotically unbiased*, if $\lim_{m \rightarrow +\infty} \text{bias}(\hat{\theta}) = 0$ where m : number of samples.

Deviation

- centered data, $x_{ic} = x_i - \bar{x}$
- average absolute deviation = $\frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$
- variance (方差) = $\frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2$
- standard deviation = $\sqrt{\frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2}$

Variance

The *variance*,

$$\begin{aligned}
 \text{Var}(X) &\doteq \frac{1}{m} \sum_{i=1}^m (X_i - E(X))^2 \\
 &= \frac{1}{m} \sum_{i=1}^m \{X_i^2 - 2E(X)X_i + [E(X)]^2\} \\
 &= \frac{1}{m} \left\{ \sum_{i=1}^m X_i^2 - 2E(X) \sum_{i=1}^m X_i + m[E(X)]^2 \right\} \\
 &= \frac{1}{m} \left\{ \sum_{i=1}^m X_i^2 - 2m[E(X)]^2 + m[E(X)]^2 \right\} \quad \left[\because E(X) = \frac{1}{m} \sum_{i=1}^m X_i \right] \\
 &= \frac{1}{m} \sum_{i=1}^m X_i^2 - [E(X)]^2 \\
 &= E(X^2) - [E(X)]^2 \\
 \Rightarrow \text{Var}(X) &= E(X^2) - [E(X)]^2 \quad \text{or} \quad \boxed{\text{Var}(X) = \langle X^2 \rangle - \langle X \rangle^2}
 \end{aligned}$$

Variance

Properties:

(1) $E(aX + b) = aE(X) + b$

Proof:

$$E(aX + b) = \sum_{i=1}^m (ax_i + b)P(x_i) = a \sum_{i=1}^m x_i p(x_i) + b \sum_{i=1}^m p(x_i) = aE(X) + b$$

□

(2) $Var(aX + b) = a^2 Var(X)$

Proof:

$$\begin{aligned} Var(aX + b) &= E((aX + b)^2) - [E(aX + b)]^2 \\ &= E(a^2 X^2 + 2abX + b^2) - [aE(X) + b]^2 \\ &= a^2 E(X^2) + 2abE(X) + b^2 - a^2 [E(X)]^2 - 2abE(X) - b^2 \\ &= a^2 [E(X^2) - [E(X)]^2] \\ &= a^2 Var(X) \end{aligned}$$

□

Sampling

Let μ : population mean

σ : population standard deviation

\bar{x} : sample mean

s : sample standard deviation

Suppose X_i ($i = 1, 2, \dots, m$) is a random sample from a population with mean μ , i.e. $E(X_i) = \mu$.

Then $E(\bar{X}) = E\left(\frac{1}{m} \sum_{i=1}^m X_i\right) = \frac{1}{m} \sum_{i=1}^m E(X_i) = \frac{1}{m} (m\mu) = \mu$.

Hence the *expectation of the sample mean* is $E(\bar{X}) = \mu$.

Sampling

$$\begin{aligned}
 \text{Var}(\bar{X}) &= \text{Var}\left(\frac{1}{m} \sum_{i=1}^m X_i\right) \\
 &= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(X_i) \\
 &= \frac{1}{m^2} (m\sigma^2) \\
 &= \frac{\sigma^2}{m}
 \end{aligned}$$

Hence the *population variance of the sample mean* is $\text{Var}(\bar{X}) = \frac{\sigma^2}{m}$.

Sampling

The *population variance*,

$$\begin{aligned}
 \sigma^2 &= \frac{1}{m} \sum_{i=1}^m (X_i - \mu)^2 \\
 &= \frac{1}{m} \sum_{i=1}^m (X_i^2 - 2\mu X_i + \mu^2) \\
 &= \frac{1}{m} \left(\sum_{i=1}^m X_i^2 - 2\mu \sum_{i=1}^m X_i + m\mu^2 \right) \\
 &= \frac{1}{m} \left(\sum_{i=1}^m X_i^2 - 2m\mu^2 + m\mu^2 \right) \\
 &= \frac{1}{m} \sum_{i=1}^m X_i^2 - \mu^2 \\
 \Rightarrow \sigma^2 &= E(X^2) - \mu^2
 \end{aligned}$$

Sampling

The (biased) sample variance is

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

Sampling

Taking the expectation on $S^2 = E(X^2) - \bar{X}^2$, the estimate (biased) sample variance is

$$\begin{aligned}
 E(S^2) &= E(E(X^2) - \bar{X}^2) \\
 &= E(X^2) - E(\bar{X}^2) \\
 &= E(X^2) - [E(\bar{X})]^2 - \left\{ E(\bar{X}^2) - [E(\bar{X})]^2 \right\} \\
 &= E(X^2) - \mu^2 - \text{Var}(\bar{X}) \\
 &= \sigma^2 - \frac{\sigma^2}{m} \\
 &= \frac{(m-1)\sigma^2}{m}
 \end{aligned}$$

Sampling

Let s be the unbiased sample variance, s.t. $E(s^2) = s^2 = \sigma^2$.

Hence $s^2 = \frac{m}{m-1} \sigma^2$

$$\Rightarrow s^2 = \frac{1}{m-1} \sum_{i=1}^m (X_i - \mu)^2$$