English

☆ Star 257

Statistics

Statistics Parameter estimation Definitions Mean estimation Variance estimation

Confidence intervals Mean Paired sample Median Trend Hypothesis testing

Paired sample Median Trend Regression analysis Least-squares estimation Correlation analysis

Probability

By Afshine Amidi and Shervine Amidi

Parameter estimation

 \square Random sample — A random sample is a collection of n random variables $X_1,...,X_n$ that are independent and identically distributed with X. □ **Estimator** — An estimator is a function of the data that is used to infer the value of an unknown parameter in a statistical model.

Blog

 \Box **Bias** — The bias of an estimator $\hat{ heta}$ is defined as being the difference between the expected value of the distribution of $\hat{ heta}$ and the true value, i.e.:

Estimating the mean

variance σ^2 , then we have:

that:

noted \overline{X} and is defined as follows:

Remark: the sample mean is unbiased, i.e $E[\overline{X}]=\mu$.

 \Box Central Limit Theorem — Let us have a random sample $X_1,...,X_n$ following a given distribution with mean μ and

$$\overline{X} \mathop{\sim}\limits_{n o +\infty} \mathcal{N}\left(\mu, rac{\sigma}{\sqrt{n}}
ight)$$

Estimating the variance

random sample is used to
$$\epsilon$$
:
$$\frac{1}{2} = \frac{n}{n} = \frac{1}{n}$$

Remark: the sample variance is unbiased, i.e $E[s^2] = \sigma^2$. \Box Chi-Squared relation with sample variance — Let s^2 be the sample variance of a random sample. We have:

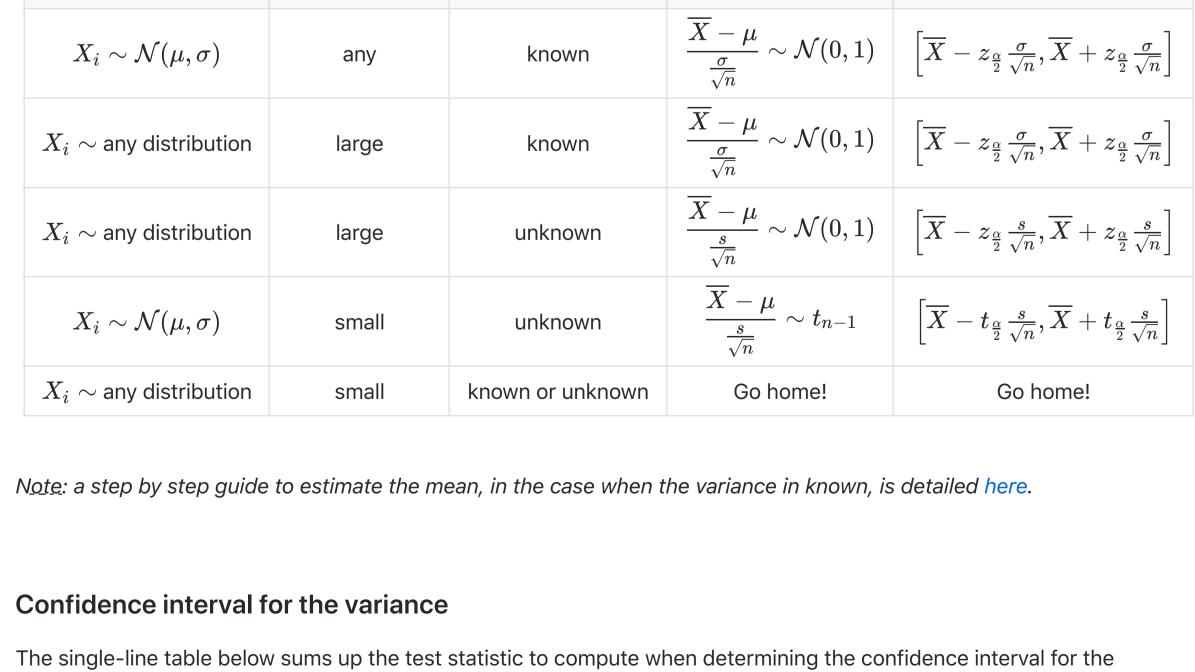
$$\left[rac{s^2(n-1)}{\sigma^2} \sim \chi^2_{n-1}
ight]$$

contained in the confidence interval.

Definitions

 \Box Confidence interval — A confidence interval $CI_{1-\alpha}$ with confidence level $1-\alpha$ of a true parameter θ is such that:

With the notation of the example above, a possible $1-\alpha$ confidence interval for θ is given by $CI_{1-\alpha}=[x_1,x_2]$.



 \Box **Type II error** — In a hypothesis test, the type II error, often noted β and also called "missed alarm", is the probability of not rejecting the null hypothesis while the null hypothesis is not true. If we note T the test statistic and R the rejection region, then we have:

rejection region, then we have:

General definitions

variance.

(left-sided)

 \Box **p-value** — In a hypothesis test, the p-value is the probability under the null hypothesis of having a test statistic T at

accept; reject

(right-sided)

p-value = $P(|T| \geqslant |T_0||H_0 ext{ true})$

 H_1

p-value = $P(T \geqslant T_0|H_0 ext{ true})$

Test statistic under $H_{
m 0}$

 \Box **Type I error** — In a hypothesis test, the type I error, often noted α and also called "false alarm" or significance level, is

 $| \alpha = P(T \in R | H_0 \text{ true}) |$

the probability of rejecting the null hypothesis while the null hypothesis is true. If we note T the test statistic and R the

$$p-$$
value

The table below sums up the test statistic to compute when performing a hypothesis test where the null hypothesis is: H_0 : $\mu_X - \mu_Y = \delta$

Sample size n_X, n_Y

any

 \square **Median of a distribution** — We define the median m of a distribution as follows:

Testing for the difference in two means

Distribution of X_i,Y_i

Normal

Testing for the median

hypothesized median.

 χ^2 test

Trends test

sequence.

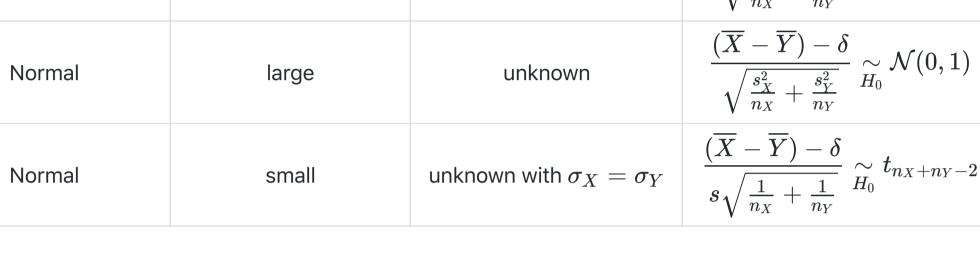
follows:

follows:

times that a larger number precedes a smaller one.

where e is referred as the error. We have:

— If np < 5, we use the following fact:



Variance σ_X^2, σ_Y^2

known

Variance σ_X^2, σ_Y^2 Distribution of X_i, Y_i Sample size $n=n_X=n_Y$ Test statistic under H_0 $rac{\overline{D} - \delta}{rac{s_D}{\sqrt{n}}} \sim_{H_0} t_{n-1}$ Normal, paired any unknown

— If $np \geqslant 5$, we use the following test statistic: $Z=rac{V-rac{n}{2}}{rac{\sqrt{n}}{2}} \mathop{\sim}\limits_{H_0} \mathcal{N}(0,1)$

 $igg| V \mathop{\sim}\limits_{H_0} \mathcal{B}\left(n,p=rac{1}{2}
ight)$

m

□ Sign test — The sign test is a non-parametric test used to determine whether the median of a sample is equal to the

By noting $V \sim \mathcal{B}(n,p=rac{1}{2})$ the number of samples falling to the right of the hypothesized median, we have:

 $\left|T = \sum_{i=1}^k rac{(Y_i - np_i)^2}{np_i} \mathop\sim_{H_0} \chi_{df}^2
ight| \quad ext{with} \quad \left[df = (k-1) - \#(ext{estimated parameters})
ight]$

Example: the sequence $\{1,5,4,3\}$ has T=3 transpositions because 5>4,5>3 and 4>3

versus

If we note x the number of transpositions in the sequence, the p-value is computed as:

determine whether the data suggest the presence of an increasing trend:

 H_0 : no trend

 \square Number of transpositions — In a given sequence, we define the number of transpositions, noted T, as the number of

□ **Test for arbitrary trends** — Given a sequence, the test for arbitrary trends is a non-parametric test, whose aim is to

p-value = $P(T \leqslant x)$

Remark: the test for a decreasing trend of a given sequence is equivalent to a test for an increasing trend of the inversed

 $Y = \alpha + \beta X + e$

 \Box Regression estimation — When estimating the regression coefficients α, β by A, B, we obtain predicted values \hat{Y}_i as

 $\hat{Y}_i = A + Bx_i$

□ **Sum of squared errors** — By keeping the same notations, we define the sum of squared errors, also known as SSE, as

 $\left| SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - (A + Bx_i))^2 \,
ight|$

Regression analysis In the following section, we will note $(x_1, Y_1), ..., (x_n, Y_n)$ a collection of n data points. \Box Simple linear model — Let X be a deterministic variable and Y a dependent random variable. In the context of a simple

 \Box Method of least-squares — The least-squares method is used to find estimates A, B of the regression coefficients α, β by minimizing the SSE. In other words, we have: $A,B = rg \min_{lpha,eta} \sum_{i=1}^n (Y_i - (lpha + eta x_i))^2 \, .$

 $S_{XY} = \sum_{i=1}^n (x_i - \overline{x})(Y_i - \overline{Y}) \quad ext{and} \quad S_{XX} = \sum_{i=1}^n (x_i - \overline{x})^2 \quad ext{and} \quad S_{YY} = \sum_{i=1}^n (Y_i - \overline{Y})^2$

 $oxed{A = \overline{Y} - rac{S_{XY}}{S_{XX}} \overline{x} \quad ext{and} \quad B = rac{S_{XY}}{S_{XX}}}$

 \Box Sum of squared errors revisited — The sum of squared errors defined above can also be written in terms of S_{YY} , S_{XY}

 $oxed{SSE = S_{YY} - BS_{XY}}$

 \Box Least-squares estimates — When estimating the coefficients α, β with the least-squares method, we obtain the

 \square **Notations** — Given n data points (x_i, Y_i) , we define S_{XY}, S_{XX} and S_{YY} as follows:

estimates A, B defined as follows:

and B as follows:

Coefficient

lpha

 β

lpha

follows:

Correlation analysis

and Y, we use the following statistic:

coefficient estimate:

View PDF version on GitHub

Sample size

Estimate

A

B

 \boldsymbol{A}

The estimator s^2 has the following property:

The table below sums up the properties surronding the least-squares estimates A,B when σ is known or not:

Statistic

known $\left| rac{A-lpha}{\sigma\sqrt{rac{1}{n}+rac{\overline{X}^2}{S_{XX}}}} \sim \mathcal{N}(0,1)
ight| \left[A-z_{rac{lpha}{2}}\sigma\sqrt{rac{1}{n}+rac{\overline{X}^2}{S_{XX}}}, A+z_{rac{lpha}{2}}\sigma\sqrt{rac{1}{n}+rac{\overline{X}^2}{S_{XX}}}
ight]$

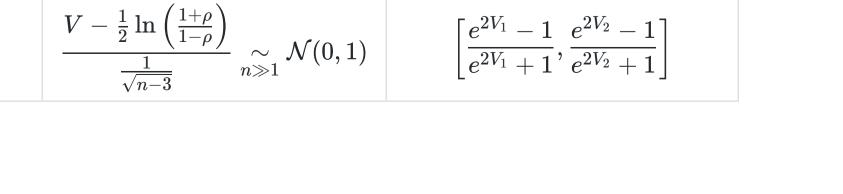
known $rac{B-eta}{rac{\sigma}{\sqrt{S_{XX}}}} \sim \mathcal{N}(0,1)$ $\left[B-z_{rac{lpha}{2}}rac{\sigma}{\sqrt{S_{XX}}},B+z_{rac{lpha}{2}}rac{\sigma}{\sqrt{S_{XX}}}
ight]$

 $ho = rac{E[(X - \mu_X)(Y - \mu_Y)]}{\sqrt{E[(X - \mu_X)^2]E[(Y - \mu_Y)^2]}}$ □ Sample correlation coefficient — The correlation coefficient is in practice estimated by the sample correlation coefficient, often noted r or $\hat{\rho}$, which is defined

By noting $V_1=V-rac{z_{rac{lpha}{2}}}{\sqrt{n-3}}$ and $V_2=V+rac{z_{rac{lpha}{2}}}{\sqrt{n-3}}$, the table below sums up the key results surrounding the correlation

□ **Fisher transformation** — The Fisher transformation is often used to build confidence intervals for correlation. It is noted

1-lpha confidence interval



1-lpha confidence interval for ho

Definitions

Key results Correlation • Fisher transformation Statistics cheatsheet

CME 106 - Introduction to Probability and Statistics for Engineers

 $ig| ext{Bias}(\hat{ heta}) = E[\hat{ heta}] - heta$ Remark: an estimator is said to be unbiased when we have $E[\hat{ heta}] = heta.$ \Box Sample mean — The sample mean of a random sample is used to estimate the true mean μ of a distribution, is often

 \Box Characteristic function for sample mean — The characteristic function for a sample mean is noted $\psi_{\overline{X}}$ and is such

 \Box Sample variance — The sample variance of a random sample is used to estimate the true variance σ^2 of a distribution, is often noted s^2 or $\hat{\sigma}^2$ and is defined as follows: $\left|s^2=\hat{\sigma}^2=rac{1}{n-1}\sum_{i=1}^n(X_i-\overline{X})^2
ight|$

Confidence intervals

 $P(heta \in CI_{1-lpha}) = 1-lpha$

Confidence interval for the mean When determining a confidence interval for the mean μ , different test statistics have to be computed depending on which case we are in. The table below sums it up. $1-\alpha$ confidence interval Distribution of X_i Sample size nVariance σ^2 **Statistic** $X_i \sim \mathcal{N}(\mu, \sigma)$ known any known $X_i \sim$ any distribution large

 $\beta = P(T \notin R|H_0 \text{ not true})$

least as extreme as the one that we observed T_0 . We have:

p-value = $P(T \leqslant T_0|H_0 ext{ true})$

Remark: the example below illustrates the case of a right-sided p-value.

(two-sided)

$$\frac{p-\text{value}}{T_0}$$

$$\square \text{ Non-parametric test} - \text{A non-parametric test is a test where we do not have any underlying assumption regarding the distribution of the sample.}$$

Testing for the mean of a paired sample We suppose here that X_i and Y_i are pairwise dependent. By noting $D_i=X_i-Y_i$, the one-line table below sums up the test statistic to compute when performing a hypothesis test where the null hypothesis is:

 H_0 : $\overline{D}=\delta$

 $P(X\leqslant m)=P(X\geqslant m)=rac{1}{2}$ $P(X \leqslant m) \triangleq \frac{1}{2}$

$$\chi^2$$
 test \square Goodness of fit test — Let us have k bins where in each of them, we observe Y_i number of samples. Our null hypothesis is that Y_i follows a binomial distribution with probability of success being p_i for each bin. We want to test whether modelling the problem as described above is reasonable given the data that we have. In order to do this, we perform a hypothesis test:
$$H_0: \operatorname{good} \operatorname{fit} \qquad \operatorname{versus} \qquad H_1: \operatorname{not} \operatorname{good} \operatorname{fit}$$

$$\square \chi^2 \operatorname{statistic} \operatorname{for} \operatorname{goodness} \operatorname{of} \operatorname{fit} - \operatorname{In} \operatorname{order} \operatorname{to} \operatorname{perform} \operatorname{the} \operatorname{goodness} \operatorname{of} \operatorname{fit} \operatorname{test}, \operatorname{we} \operatorname{need} \operatorname{to} \operatorname{compute} \operatorname{a} \operatorname{test} \operatorname{statistic} \operatorname{that}$$
 we can compare to a reference distribution. By noting k the number of bins, n the total number of samples, if we have $np_i \geqslant 5$, the test statistic T defined below will enable us to perform the hypothesis test:

linear model, we assume that Y is linked to X via the regression coefficients lpha,eta and a random variable $e\sim\mathcal{N}(0,\sigma)$,

 $|H_1|$: there is an increasing trend

Key results When
$$\sigma$$
 is unknown, this parameter is estimated by the unbiased estimator s^2 defined as follows:
$$s^2 = \frac{S_{YY} - BS_{XY}}{n-2}$$
 The estimator s^2 has the following property:

 $\begin{array}{ll} \text{unknown} & \frac{A-\alpha}{s\sqrt{\frac{1}{n}+\frac{\overline{X}^2}{S_{XX}}}}\sim t_{n-2} & \left[A-t_{\frac{\alpha}{2}}s\sqrt{\frac{1}{n}+\frac{\overline{X}^2}{S_{XX}}},A+t_{\frac{\alpha}{2}}s\sqrt{\frac{1}{n}+\frac{\overline{X}^2}{S_{XX}}}\right] \\ \text{unknown} & \frac{B-\beta}{\frac{s}{\sqrt{S_{XX}}}}\sim t_{n-2} & \left[B-t_{\frac{\alpha}{2}}\frac{s}{\sqrt{S_{XX}}},B+t_{\frac{\alpha}{2}}\frac{s}{\sqrt{S_{XX}}}\right] \end{array}$ β B

coefficient, often noted r or $\hat{ ho}$, which is defined a	is:
	$r=\hat{ ho}=rac{S_{XY}}{\sqrt{S_{XX}S_{YY}}}$
□ Testing for correlation — In order to perform	a hypothesis test with H_0 being that there is no correlation between X

 \Box Correlation coefficient — The correlation coefficient of two random variables X and Y is noted ρ and is defined as

V and defined as follows: $V = rac{1}{2} \ln \left(rac{1+r}{1-r}
ight)$

Standardized statistic

y in G G