

Statistics 1 Notes

paraphrased by Tyler Wright

*An important note, these notes are absolutely **NOT** guaranteed to be correct, representative of the course, or rigorous. Any result of this is not the author's fault.*

1 The Basics of Data Analysis

1.1 Samples

A sample is a set of values observed from a simple random sample of some size n from a population where each sample member is chosen **independently** of each other and each population member is **equally likely** to be selected.

Samples are usually written as $\{x_1, x_2, \dots, x_n\}$ where each x_i represents an observed value. If the data is ordered, the data is written as $\{x_{(1)}, x_{(2)}, \dots, x_{(n)}\}$ (for numerical values, this is ascending order). So, in this case, $x_{(1)}$ is always the minimum, $x_{(n)}$ is always the maximum.

1.2 Probability Density Functions

For a sample $\{x_1, x_2, \dots, x_n\}$, we can imagine each datum as being distributed with some population distribution X . As each datum is independent of all other observed values, we can write the probability density of this sample as follows:

$$f_X(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_X(x_i).$$

1.3 Measures of Central Tendency

1.3.1 Sample Median

For a sample $X = \{x_1, x_2, \dots, x_n\}$, we define the sample median M as follows:

$$M(X) = \begin{cases} x_{(m+1)} & \text{for } n = 2m + 1 \\ \frac{x_{(m)} + x_{(m+1)}}{2} & \text{for } n = 2m. \end{cases}$$

Essentially, it equals the middle value or the average of the middle values. Also, it's important to note that the median is not sensitive to extreme values.

1.3.2 Sample Mean

For a sample $X = \{x_1, x_2, \dots, x_n\}$, we define the sample mean \bar{X} as follows:

$$\bar{X} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right).$$

This is easy to calculate even when combining samples. However, it is sensitive to extreme values.

1.3.3 Trimmed Sample Mean

For a sample $X = \{x_1, x_2, \dots, x_n\}$, we define the trimmed sample mean \overline{X}_Δ for some percentage $\Delta\%$ as follows:

$$\begin{aligned}\text{Let } k &= \left\lfloor n \frac{\Delta}{100} \right\rfloor \\ \text{Let } \tilde{X} &= \{x_{(k+1)}, x_{(k+2)}, \dots, x_{(n-k)}\} \\ \overline{X}_\Delta &= \overline{\tilde{X}} \text{ (the sample mean of } \tilde{X}\text{).}\end{aligned}$$

Basically, you remove the first and last $\Delta\%$ of values and take the sample mean of the remaining values.

1.4 Measures of Spread

1.4.1 Sample Variance

For a sample $X = \{x_1, x_2, \dots, x_n\}$, we define the sample variance s^2 as follows:

$$\begin{aligned}s^2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right).\end{aligned}$$

This measures how much the data varies.

1.4.2 Hinges

There are two hinge measures, lower (H_1) and upper (H_3):

$$\begin{aligned}H_1 &= \text{median of \{data values } \leq \text{ the median\}} \\ H_3 &= \text{median of \{data values } \geq \text{ the median\}}.\end{aligned}$$

1.4.3 Quartiles

For a sample $X = \{x_1, x_2, \dots, x_n\}$, there are two quartile measures, lower (Q_1) and upper (Q_3). The formulas are long and overly complicated so for Q_1 :

- Calculate $k = \frac{n+1}{4}$
- If $k \in \mathbb{Z}$, $Q_1 = x_{(k)}$
- Otherwise, do linear interpolation between $x_{(\lfloor k \rfloor)}$ and $x_{(\lfloor k \rfloor + 1)}$

And similarly for Q_3 :

- Calculate $k = 3 \left(\frac{n+1}{4} \right)$
- If $k \in \mathbb{Z}$, $Q_3 = x_{(k)}$
- Otherwise, do linear interpolation between $x_{(\lfloor k \rfloor)}$ and $x_{(\lfloor k \rfloor + 1)}$

For large samples, the quartiles and hinges tend to be close to each other.

1.4.4 Interquartile Range (IQR)

The IQR is the difference between Q_3 and Q_1 ($Q_3 - Q_1$).

In this course, outliers are defined as more than $\frac{3}{2}(\text{IQR})$ (or approx. $\frac{3}{2}(H_3 - H_1)$) from the median.

1.4.5 Skewness

We measure skewness by the distance of the hinges from the median. If H_3 is further from the median than H_1 , we have a longer right tail. If the converse is true, we have a longer left tail.

2 Assessing Fit

2.1 Quantiles of a Distribution

For a distribution X with cumulative distribution function F_X , the quantiles of the distribution are defined as the set of values:

$$F_X^{-1} \left\{ \frac{1}{n+1}, \dots, \frac{n}{n+1} \right\}.$$

We use $n+1$ on the denominator as $F_X^{-1}(1)$ can be ∞ .

The ordered sample is called the set of sample quantiles.

2.2 Quantile-Quantile (Q-Q) Plots

These are the steps for constructing a Q-Q plot of a sample $\{x_1, x_2, \dots, x_n\}$ with cumulative distribution function F_X :

- Generate an estimate for the parameter(s) $(\hat{\theta}_1, \hat{\theta}_2, \dots)$
- Compute the quantiles (the expected quantiles if the hypothesised model is correct)
- Plot each expected quantile against the sample quantile $(F_X^{-1}(\frac{k}{n+1}; \hat{\theta}), x_{(k)})$.

What we would expect, if our hypothesis is correct, is that the plotted points lie close to the line $y = x$. This is saying our sample and expected quantiles are close together.

2.3 Probability Plots

These are similar to the Q-Q plots but plot the sample cumulative probability against expected probability $(F_X(x_{(k)}), \frac{k}{n+1})$.

3 Estimation

We have that a population is distributed with some distribution X with a probability density function (PDF) f_X , cumulative distribution function (CDF) F_X , and some parameters $\{\theta_1, \dots\}$. We can make guesses at the distribution of a sample and use tests to verify that. But, to do these tests we need a value for the parameters. It's not practical to guess these, so we need to estimate them.

3.1 Parameters

We say $\hat{\theta}$ is an estimator for θ and define it as a function of a sample $\{x_1, x_2, \dots, x_n\}$:

$$\hat{\theta}(x_1, x_2, \dots, x_n).$$

3.2 Distribution Quantities

From our estimated value of the distribution parameters, we can calculate estimated values for distribution quantities like the mean and variance. We consider τ a function of the parameter that gives a distribution quantity:

- **True quantity:** $\tau(\theta)$ where θ is the true distribution parameter
- **Estimated quantity:** $\hat{\tau} = \tau(\hat{\theta})$ where $\hat{\theta}$ is our estimated parameter.

4 Method of Moments Estimation

4.1 Definition of a Moment

The k th moment of a probability distribution X is defined as follows:

$$\mathbb{E}(X^k) := \int_{-\infty}^{\infty} x^k f_X(x) dx.$$

Setting $k = 1$ gives us the familiar expectation of X :

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f(x) dx.$$

In the discrete case, the integral is a sum.

We define the k th sample moment m_k as follows:

$$m_k = \frac{\sum_{i=1}^n x_i^k}{n}.$$

Or rather, the k th moment is the average value of x^k in the sample.

4.2 The Process

By considering the a probability distribution X with parameter θ , we can find functions for the moments of X in terms of θ . These can be rearranged to give functions for θ in terms of the moments. We can then use the sample moments to generate an estimate for θ ($\hat{\theta}_{mom}$).

4.3 Method of Moments on the Exponential

Assume we have some population X distributed according to the Exponential with some parameter θ . We say $X \sim \text{Exp}(\theta)$.

$$\begin{aligned} f_X(x) &= \theta e^{-\theta x} && (x > 0) \\ \Rightarrow \mathbb{E}(X) &= \frac{1}{\theta} \\ \Rightarrow \theta &= \frac{1}{\mathbb{E}(X)} \\ \Rightarrow \hat{\theta}_{mom} &= \frac{1}{m_1}. \end{aligned}$$

If there were more parameters, we would have to consider greater moments of X .

5 Maximum Likelihood Estimation

5.1 The Process

By considering the a probability distribution X with parameter θ , we can find functions for the probability of events occuring in terms of θ . If we find where this function is maximised, it will give us the value of θ that makes this sample most likely. This is the maximum likelihood estimate ($\hat{\theta}_{mle}$).

5.2 Optimisation of the Method

Consider a sample $\{x_1, x_2, \dots, x_n\}$ with distributions $\{X_1, X_2, \dots, X_n\}$, we call the likelihood function the joint PDF of $\{X_1, X_2, \dots, X_n\}$. We input our sample values ($L = f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n)$) which gives us a function in terms of our unknown parameters $\theta_1, \theta_2, \dots$.

For X_1, X_2, \dots, X_n independent and identically distributed sharing some distribution X , the joint PDF can be written as a product of marginals:

$$\begin{aligned} L &= f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) \\ &= f_X(x_1)f_X(x_2) \cdots f_X(x_n) \\ &= \prod_{i=1}^n f_X(x_i). \end{aligned}$$

We can take the natural logarithm of this likelihood function (the value where it's maximised is preserved as the natural logarithm is increasing) ($\ell = \ln(L)$). If, again, X_1, X_2, \dots, X_n are independent and identically distributed sharing some distribution X :

$$\begin{aligned} \ell &= \ln \left(\prod_{i=1}^n f_X(x_i) \right) \\ &= \sum_{i=1}^n [\ln (f_X(x_i))] \end{aligned}$$

We know $\hat{\theta}_{mle}$ is the solution to:

$$\frac{\partial}{\partial \theta} \ell(\theta) = 0.$$

So, in the independent and identically distributed case:

$$\frac{\partial}{\partial \theta} \ell(\theta) = 0 \iff \sum_{i=1}^n \left[\frac{\partial}{\partial \theta} \ln (f_X(x_i)) \right] = 0.$$

5.3 Multiple Parameters

When finding the maximum likelihood estimate for multiple parameters, we obtain multiple equations by partially differentiating ℓ by our different parameters, giving an equation for each.

5.4 Non-regular density

If our function L is piecewise, we may find that our maximum isn't where the derivative is zero, but as the endpoints of the parts of the function.

Consider the maximum of $f : \mathbb{R} \rightarrow \mathbb{R}$ where:

$$f(x) = \begin{cases} \theta^{-x} & \text{for } x \geq 1 \\ 0 & \text{otherwise.} \end{cases}$$

6 The Performance of Estimators

6.1 Variation of Estimators

We can consider the distribution of an estimator to compare them. We consider these main quantities:

- $\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$
- $\text{mse}(\hat{\theta}) = \mathbb{E}((\hat{\theta} - \theta)^2)$.

Bias is the quantity we expect our estimator to vary by from the true value. The MSE (mean squared error) is how much it varies.

We can rewrite the formula for the mean squared error as follows:

$$\text{mse}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{bias}(\hat{\theta})^2.$$

6.2 Method of Simulation

If we have a distribution with known parameters $\theta_1, \theta_2, \dots$, we can sample N samples of size n and use our estimators to calculate estimates for these known parameters for each sample.

From this, we can calculate the average error, sample variance, and average squared error of each estimator. These quantities estimate bias, variance, and mean squared error respectively.

If we repeat this process for multiple estimators, we can compare our estimators with these quantities.

7 Central Limit Theorem

7.1 Definition of the Central Limit Theorem

For X_1, X_2, \dots, X_n a random sample from a population with mean $\mu = \mathbb{E}(X)$ and variance $\sigma^2 = \text{Var}(X)$. Let \overline{X}_n be the sample mean. For n large we have:

$$\mathbb{P}\left(\sqrt{n} \left[\frac{\overline{X}_n - \mu}{\sigma} \right] \leq x\right) \simeq \mathbb{P}(\mathcal{N}(0, 1) \leq x) = \Phi(x).$$

Or similarly:

$$\overline{X}_n \simeq \mathcal{N}(\mu, \sigma^2/n).$$

7.2 Continuity Correction

When using the Central Limit Theorem to approximate discrete random variables, it is important to make a continuity correction. Let X_1, X_2, \dots, X_n be samples from a discrete random variable with sample mean \overline{X}_n , population mean μ , and population variance σ^2 :

$$\begin{aligned} \mathbb{P}\left(\sqrt{n} \left[\frac{\overline{X}_n - \mu}{\sigma} \right] \leq x\right) &\simeq \mathbb{P}(\mathcal{N}(0, 1) < x + \frac{1}{2}) \\ \mathbb{P}\left(\sqrt{n} \left[\frac{\overline{X}_n - \mu}{\sigma} \right] < x\right) &\simeq \mathbb{P}(\mathcal{N}(0, 1) < x - \frac{1}{2}) \end{aligned}$$

8 A Reminder on Moment Generating Functions

8.1 Definition of a Moment Generating Function (MGF)

For a random variable X , we define the moment generating function by:

$$\mathcal{M}_X(t) := \mathbb{E}(e^{tX}) = \begin{cases} \int_{-\infty}^{\infty} e^{tX} f_X(x) dx & \text{for } X \text{ continuous} \\ \sum_{x \in S} e^{tX} \mathbb{P}(X = x) & \text{for } X \text{ discrete.} \end{cases}$$

8.2 Properties of a Moment Generating Function

8.2.1 Standard examples of moment generating functions

For a random variable X :

- $X \sim \mathcal{N}(\mu, \sigma^2) \Leftrightarrow \mathcal{M}_X(t) = \exp(\mu t + \frac{(\sigma t)^2}{2})$
- $X \sim \text{Exp}(\theta) \Leftrightarrow \mathcal{M}_X(t) = \frac{\theta}{\theta - t}$
- $X \sim \text{Gamma}(\alpha, \beta) \Leftrightarrow \mathcal{M}_X(t) = \frac{\beta^\alpha}{(\beta - t)^\alpha}$.

8.2.2 Joint moment generating functions

The joint MGF of X and Y is:

$$\mathcal{M}_{X,Y}(s, t) := \mathbb{E}(e^{sX+tY}).$$

They are such that:

$$\begin{aligned}\mathcal{M}_X(s) &= \mathcal{M}_{X,Y}(s, 0) \\ \mathcal{M}_Y(t) &= \mathcal{M}_{X,Y}(0, t).\end{aligned}$$

We also have that X and Y are independent if and only if:

$$\mathcal{M}_{X,Y}(s, t) = \mathcal{M}_X(s)\mathcal{M}_Y(t).$$

8.2.3 Independence of moment generating functions

If X_1, X_2, \dots, X_n are independent and $Y = \sum_{i=1}^n X_i$:

$$\mathcal{M}_Y(t) = \prod_{i=1}^n \mathcal{M}_{X_i}(t)$$

8.2.4 Uniqueness of moment generating functions

The MGF uniquely defines a distribution, for two random variables X, Y :

$$\mathcal{M}_X = \mathcal{M}_Y \Leftrightarrow X = Y.$$

9 The Normal Distribution

9.1 Transformation and Addition of the Normal

For $X \sim \mathcal{N}(\mu, \sigma^2)$ and $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ for $i \in \{1, 2, \dots, n\}$, let $\bar{X} = \frac{1}{n}(\sum_{i=1}^n X_i)$ be the sample mean:

$$\begin{aligned}
 aX + b &\sim \mathcal{N}(a\mu + b, a^2\sigma^2) && \text{(Linear Transformation)} \\
 \sum_{i=1}^n X_i &\sim \mathcal{N}\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right) && \text{(Summed)} \\
 \frac{(X - \mu)}{\sigma} &\sim \mathcal{N}(0, 1) && \text{(Standardised)} \\
 \bar{X} &\sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) && \text{(Sample Mean)} \\
 \sqrt{n} \left(\frac{\bar{X} - \mu}{\sigma} \right) &\sim \mathcal{N}(0, 1). && \text{(Standardised Sample Mean)}
 \end{aligned}$$

It's very important to remember that multiplication and summing differ when dealing with the Normal (when it comes to the variance). So, if you have a Normal random variable $X \sim \mathcal{N}(\mu, \sigma^2)$, $2X \neq X + X$ as $\text{Var}(2X) = 4\sigma^2$ and $\text{Var}(X + X) = 2\sigma^2$. This is because when you're multiplying, you're amplifying variation in your sample, but when you sum you're combining variance.

9.2 Independence of the Sample Mean and the Sum of Squared Difference

For $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ for $i \in \{1, 2, \dots, n\}$, let $\bar{X} = \frac{1}{n}(\sum_{i=1}^n X_i)$ be the sample mean. We have that \bar{X} and $\sum_{i=1}^n (X_i - \bar{X})^2$ are independent.

10 Sampling Distributions related to the Normal

10.1 The χ^2 Distribution

10.1.1 Definition of the χ^2 distribution

We say that a random variable $X \sim \chi_r^2$ (r degrees of freedom) if:

$$\mathcal{M}_X(t) = (1 - 2t)^{-r/2}.$$

10.1.2 Properties of the χ^2 distribution

Let $X \sim \chi_r^2$, $Y \sim \chi_s^2$:

- $X \sim \Gamma(\frac{r}{2}, \frac{1}{2})$
- $\mathbb{E}(X) = r$
- $\text{Var}(X) = 2r$
- $X + Y \sim \chi_{r+s}^2$

We also have some results relating to the Normal, let Z be the standard Normal, X_i for $i \in \{1, 2, \dots, n\}$ be samples from $\mathcal{N}(\mu, \sigma^2)$, let $\bar{X} = \frac{1}{n}(\sum_{i=1}^n X_i)$ be the sample mean:

- $Z^2 \sim \chi_1^2$
- $\sum_{i=1}^n (\frac{X_i - \mu}{\sigma})^2 \sim \chi_n^2$
- $\sum_{i=1}^n (\frac{X_i - \bar{X}}{\sigma})^2 \sim \chi_{n-1}^2$.

Finally, we also have some results relating to the Exponential and Gamma distributions, let X_i for $i \in \{1, 2, \dots, n\}$ be samples from $\text{Exp}(\theta)$:

- $\sum_{i=1}^n X_i \sim \Gamma(n, \theta)$
- $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim \Gamma(n, n\theta)$
- $2\theta \sum_{i=1}^n X_i \sim \Gamma(n, 1/2) = \chi_{2n}^2$

10.2 The t Distribution

10.2.1 Definition of the t distribution

For $Z \sim \mathcal{N}(0, 1)$, $X \sim \chi_r^2$ **independent** we have:

$$T = \frac{Z}{\sqrt{X/r}},$$

is distributed with a t distribution with r degrees of freedom (t_r).

10.2.2 Properties of the t distribution

For $T \sim t_r$:

- $\mathbb{E}(T) = 0$
- $\text{Var}(T) = \frac{r}{r-2}$
- The density of T approaches $\mathcal{N}(0, 1)$ as $r \rightarrow \infty$.

10.2.3 Samples from the Normal with σ unknown

For X_i for $i \in \{1, 2, \dots, n\}$ be samples from $\mathcal{N}(\mu, \sigma^2)$, let $\bar{X} = \frac{1}{n}(\sum_{i=1}^n X_i)$ be the sample mean and $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ be the sample variance. We have that:

$$\sqrt{n} \left(\frac{\bar{X} - \mu}{S} \right) \sim t_{n-1}$$

*This is **extremely** key as this allows us to perform hypothesis test on any Normal sample without knowing the population variance.*

11 Confidence Intervals

11.1 Definition of a Confidence Interval

For $X_1, \dots, X_n \sim f_{X_1, \dots, X_n}$ with some unknown parameter θ , we have the $\alpha\%$ confidence interval for θ is (c_L, c_U) where:

$$\mathbb{P}(c_L \leq \theta \leq c_U) \geq \alpha.$$

So, $\alpha\%$ of intervals constructed in this way contain the true value of θ .

11.2 Examples of Confidence Intervals

11.2.1 $\mathcal{N}(\mu, \sigma^2)$: Confidence interval for μ with σ^2 known

The $\alpha\%$ confidence interval for μ from a sample $\{x_1, x_2, \dots, x_n\}$ is (c_L, c_U) where Φ^{-1} is the Normal quantile function (inverse CDF) and:

$$\begin{aligned} c_L &= \bar{x} - \Phi^{-1} \left(\frac{1 + \alpha}{2} \right) \left(\frac{\sigma}{\sqrt{n}} \right) \\ c_U &= \bar{x} + \Phi^{-1} \left(\frac{1 + \alpha}{2} \right) \left(\frac{\sigma}{\sqrt{n}} \right). \end{aligned}$$

So, for a 95% confidence interval:

$$\begin{aligned} c_L &= \bar{x} - \Phi^{-1}(0.975) \left(\frac{\sigma}{\sqrt{n}} \right) \\ c_U &= \bar{x} + \Phi^{-1}(0.975) \left(\frac{\sigma}{\sqrt{n}} \right). \end{aligned}$$

11.2.2 $\mathcal{N}(\mu, \sigma^2)$: Confidence interval for μ with σ^2 unknown

The $\alpha\%$ confidence interval for μ from a sample $\{x_1, x_2, \dots, x_n\}$ is (c_L, c_U) where $\hat{\sigma}^2$ is the sample variance, T_n^{-1} is the t_n quantile function (inverse CDF), and:

$$\begin{aligned} c_L &= \bar{x} - T_{n-1}^{-1} \left(\frac{1 + \alpha}{2} \right) \left(\frac{\hat{\sigma}}{\sqrt{n}} \right) \\ c_U &= \bar{x} + T_{n-1}^{-1} \left(\frac{1 + \alpha}{2} \right) \left(\frac{\hat{\sigma}}{\sqrt{n}} \right). \end{aligned}$$

11.2.3 $\mathcal{N}(\mu, \sigma^2)$: Confidence interval for σ^2 with μ unknown

The $\alpha\%$ confidence interval for σ^2 from a sample $\{x_1, x_2, \dots, x_n\}$ is (c_L, c_U) where X_n^{-1} is the χ_n^2 quantile function (inverse CDF):

$$c_L = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{X_{n-1}^{-1}\left(\frac{1-\alpha}{2}\right)}$$
$$c_U = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{X_{n-1}^{-1}\left(\frac{1+\alpha}{2}\right)}.$$

11.2.4 $U(0, \theta)$: Confidence interval for θ

The $\alpha\%$ confidence interval for θ from a sample $\{x_1, x_2, \dots, x_n\}$ is (c_L, c_U) where:

$$c_L = x_{(n)} \left(\frac{1+\alpha}{2} \right)^{-\frac{1}{n}}$$
$$c_U = x_{(n)} \left(\frac{1-\alpha}{2} \right)^{-\frac{1}{n}}.$$

11.2.5 $\text{Exp}(\theta)$: Confidence interval for θ

The $\alpha\%$ confidence interval for θ from a sample $\{x_1, x_2, \dots, x_n\}$ is (c_L, c_U) where X_n^{-1} is the χ_n^2 quantile function (inverse CDF):

$$c_L = \frac{X_{2n}^{-1}\left(\frac{1+\alpha}{2}\right)}{2 \sum_{i=1}^n x_i}$$
$$c_U = \frac{X_{2n}^{-1}\left(\frac{1-\alpha}{2}\right)}{2 \sum_{i=1}^n x_i}.$$

11.2.6 Confidence Intervals by Simulation

For a distribution X with unknown parameter θ , we can find a $\alpha\%$ confidence interval (c_L, c_U) by simulation from a sample $\{x_1, x_2, \dots, x_n\}$:

- Calculate an estimate for θ ($\hat{\theta}$) by using the sample
- Use the estimate as a parameter and generate N samples of some size
- Calculate the estimates for θ these samples
- Find the values c_L and c_U such that $\alpha\%$ of the estimates are in (c_L, c_U) .

12 Notes on Hypothesis Testing: Population Means

12.1 Test Statistics

Test statistics are generated with the assumption of H_0 being true. This means if it's unlikely the observed value is distributed as expected, we can reject our H_0 in favour of an alternative hypothesis. This is the basis of our hypothesis tests.

12.2 p -values

p -values are calculated from an observed test statistic t_{obs} by considering the probability that the distribution of our test statistic T is further from what is expected than our observed value:

$$p\text{-value} = \begin{cases} \mathbb{P}(|T| \geq |t_{obs}|) & (H_1 : \neq) \\ \mathbb{P}(T \geq t_{obs}) & (H_1 : >) \\ \mathbb{P}(T \leq t_{obs}) & (H_1 : <). \end{cases}$$

Small p -values indicate that it's unlikely that our test statistic would be further from what is expected than t_{obs} . This means our assumption of H_0 may not be accurate.

12.3 Error

There's two main types of error in hypothesis testing:

- **Type I error:** The probability that we conclude H_0 is false when it is in fact true
- **Type II error:** The probability that we conclude H_0 is true when H_1 is in fact true.

We can fix our Type I error to a value α before conducting our test by using a critical region with a $\alpha\%$ significance level.

We call the power of a test P where:

$$P = 1 - (\text{Type II error}).$$

13 Hypothesis Tests: Comparison of Population Means

13.1 Normal Hypothesis Test

13.1.1 Assumptions

For a sample $\{x_1, \dots, x_n\}$ randomly sampled from $\mathcal{N}(\mu, \sigma^2)$ with σ^2 known and an expected mean μ_0 .

13.1.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \mu &= \mu_0 \\ H_1 : \mu &\neq \text{or } > \text{ or } < \mu_0. \end{aligned}$$

13.1.3 Test Statistic

We have the test statistic defined by:

$$T = \frac{\sqrt{n}(\bar{x} - \mu_0)}{\sigma} \sim \mathcal{N}(0, 1).$$

13.1.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let Φ^{-1} be the $\mathcal{N}(0, 1)$ quantile function (inverse CDF), α be our significance level:

$$\begin{aligned} H_1 : \mu \neq \mu_0 &\Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = \Phi^{-1}(1 - \frac{\alpha}{2})) \\ H_1 : \mu > \mu_0 &\Rightarrow C = [c^*, \infty) & (c^* = \Phi^{-1}(1 - \alpha)) \\ H_1 : \mu < \mu_0 &\Rightarrow C = (-\infty, -c^*]. & (c^* = \Phi^{-1}(1 - \alpha)) \end{aligned}$$

13.2 One Sample t -test

13.2.1 Assumptions

For a sample $\{x_1, \dots, x_n\}$ randomly sampled from $\mathcal{N}(\mu, \sigma^2)$ with σ^2 unknown and an expected mean μ_0 .

13.2.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \mu &= \mu_0 \\ H_1 : \mu &\neq \text{ or } > \text{ or } < \mu_0. \end{aligned}$$

13.2.3 Test Statistic

For S^2 defined as the sample variance of our sample, we have the test statistic defined by:

$$T = \frac{\sqrt{n}(\bar{x} - \mu_0)}{S} \sim t_{n-1}.$$

13.2.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let T_k^{-1} be the t_k quantile function (inverse CDF), α be our significance level:

$$\begin{aligned} H_1 : \mu \neq \mu_0 &\Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = T_{n-1}^{-1}(1 - \frac{\alpha}{2})) \\ H_1 : \mu > \mu_0 &\Rightarrow C = [c^*, \infty) & (c^* = T_{n-1}^{-1}(1 - \alpha)) \\ H_1 : \mu < \mu_0 &\Rightarrow C = (-\infty, -c^*]. & (c^* = T_{n-1}^{-1}(1 - \alpha)) \end{aligned}$$

13.3 Pooled Two Sample t -test

13.3.1 Assumptions

For two **independent** samples with the **same variance**:

- $\{x_1, \dots, x_n\}$ randomly sampled from $\mathcal{N}(\mu_X, \sigma^2)$
- $\{y_1, \dots, y_m\}$ randomly sampled from $\mathcal{N}(\mu_Y, \sigma^2)$.

13.3.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \mu_X - \mu_Y &= 0 \text{ (the values are the same)} \\ H_1 : \mu_X - \mu_Y &\neq 0 \text{ (the values are not the same)} \\ &> 0 \text{ } (\mu_X \text{ is greater)} \\ &< 0 \text{ } (\mu_Y \text{ is greater}). \end{aligned}$$

13.3.3 Test Statistic

We define the pooled variance S_p :

$$S_p = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^m (y_i - \bar{y})^2}{n + m - 2} = \frac{S_{xx} + S_{yy}}{n + m - 2},$$

and we have the test statistic T defined by:

$$T = \frac{\bar{x} - \bar{y}}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}.$$

13.3.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let T_k^{-1} be the t_k quantile function (inverse CDF), α be our significance level:

$$\begin{aligned} H_1 : \mu_X - \mu_Y &\neq 0 \Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = T_{n+m-2}^{-1}(1 - \frac{\alpha}{2})) \\ H_1 : \mu_X - \mu_Y &> 0 \Rightarrow C = [c^*, \infty) & (c^* = T_{n+m-2}^{-1}(1 - \alpha)) \\ H_1 : \mu_X - \mu_Y &< 0 \Rightarrow C = (-\infty, -c^*]. & (c^* = T_{n+m-2}^{-1}(1 - \alpha)) \end{aligned}$$

13.4 Welch Two Sample t -test

13.4.1 Assumptions

For two **independent** samples (if the variances are the same, use the pooled test):

- $\{x_1, \dots, x_n\}$ randomly sampled from $\mathcal{N}(\mu_X, \sigma^2)$
- $\{y_1, \dots, y_m\}$ randomly sampled from $\mathcal{N}(\mu_Y, \sigma^2)$.

13.4.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \mu_X - \mu_Y &= 0 \text{ (the values are the same)} \\ H_1 : \mu_X - \mu_Y &\neq 0 \text{ (the values are not the same)} \\ &> 0 \text{ } (\mu_X \text{ is greater)} \\ &< 0 \text{ } (\mu_Y \text{ is greater}). \end{aligned}$$

13.4.3 Test Statistic

Where S_X^2, S_Y^2 are the sample variances of X and Y respectively, we have the test statistic T defined by:

$$T = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}} \simeq t_\nu,$$

where ν is defined as follows:

$$\nu = \frac{(\frac{S_X^2}{n} + \frac{S_Y^2}{m})^2}{\frac{1}{n-1}(\frac{S_X^2}{n})^2 + \frac{1}{m-1}(\frac{S_Y^2}{m})^2}$$

13.4.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let T_k^{-1} be the t_k quantile function (inverse CDF), α be our significance level:

$$\begin{aligned} H_1 : \mu_X - \mu_Y &\neq 0 \Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = T_\nu^{-1}(1 - \frac{\alpha}{2})) \\ H_1 : \mu_X - \mu_Y &> 0 \Rightarrow C = [c^*, \infty) & (c^* = T_\nu^{-1}(1 - \alpha)) \\ H_1 : \mu_X - \mu_Y &< 0 \Rightarrow C = (-\infty, -c^*] & (c^* = T_\nu^{-1}(1 - \alpha)) \end{aligned}$$

13.5 Paired t -test

13.5.1 Assumptions

For two samples, $\{x_1, \dots, x_n\}$ and $\{y_1, \dots, y_n\}$, both randomly sampled, we create the sample:

$$\{w_1, \dots, w_n\} \text{ where } w_i = x_i - y_i. \quad (i \in \{1, \dots, n\})$$

We assume each w_i is drawn from a $\mathcal{N}(\delta, \sigma^2)$ distribution where δ and σ are unknown.

13.5.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \delta &= 0 \text{ (there's no underlying difference)} \\ H_1 : \delta &\neq 0 \text{ (there's some underlying difference)} \\ &> 0 \text{ (} X \text{ tends to be greater than } Y \text{)} \\ &< 0 \text{ (} Y \text{ tends to be greater than } X \text{)}. \end{aligned}$$

13.5.3 Test Statistic

Where S_W^2 , is the sample variances of W , we have the test statistic T defined by:

$$T = \frac{\sqrt{n} \bar{w}}{S_W} \sim t_{n-1}.$$

13.5.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let T_k^{-1} be the t_k quantile function (inverse CDF), α be our significance level:

$$\begin{aligned} H_1 : \delta &\neq 0 \Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = T_{n-1}^{-1}(1 - \frac{\alpha}{2})) \\ H_1 : \delta &> 0 \Rightarrow C = [c^*, \infty) & (c^* = T_{n-1}^{-1}(1 - \alpha)) \\ H_1 : \delta &< 0 \Rightarrow C = (-\infty, -c^*]. & (c^* = T_{n-1}^{-1}(1 - \alpha)) \end{aligned}$$

14 Linear Regression

14.1 Model Assumptions

For $\{x_1, \dots, x_n\}$ values of observed of a predictor variable X with $y_i (i \in \{1, \dots, n\})$ the values of the response variable Y :

$$y_i = \alpha + \beta x_i + e_i,$$

where α and β are unknown and each e_i represents the error.

14.2 Errors

For the error values $\{e_1, \dots, e_n\}$, we assume that:

- $\mathbb{E}(e_i) = 0$
- $\text{Var}(e_i) = \sigma^2$ (unknown)
- $\text{Cov}(e_i, e_j) = 0$ for $i \neq j$ (uncorrelated error).

As a consequence, we have that:

- $\mathbb{E}(y_i) = \alpha + \beta x_i$
- $\text{Var}(y_i) = \sigma^2$
- $\text{Cov}(y_i, y_j) = 0$ for $i \neq j$ (uncorrelated responses).

14.3 Summary Statistics

We define a few summary statistics used in our calculations:

- $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$
- $S_{xx} = \sum_{i=1}^n (x_i)^2 - n\bar{x}^2$
- $S_{xy} = \sum_{i=1}^n (x_i y_i) - n\bar{x}\bar{y}$
- $S_{yy} = \sum_{i=1}^n (y_i)^2 - n\bar{y}^2$.

14.4 Finding $\hat{\alpha}$ and $\hat{\beta}$

By minimising the squared error between our estimated and observed values, we develop least squares estimates for α and β :

- $\hat{\beta} = \frac{S_{xy}}{S_{xx}}$
- $\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$

14.5 Residual Sum of Squares and $\hat{\sigma}^2$

We define the residual sum of squares RSS as the sum of all the errors squared, an alternate formula for this is:

$$RSS = S_{yy} - \frac{S_{xy}^2}{S_{xx}}.$$

Then, we can estimate σ^2 by:

$$\hat{\sigma}^2 = \frac{RSS}{n-2} = \frac{S_{yy} - \frac{S_{xy}^2}{S_{xx}}}{n-2}.$$

15 Notes on Hypothesis Testing: Regression

15.1 Assumption of Normality

To perform linear regression hypothesis tests, we use an assumption of normality. We assume the errors e_i are independent and identically distributed $\mathcal{N}(0, \sigma^2)$.

15.2 The Distribution of $\hat{\alpha}$, $\hat{\beta}$, and $\hat{\sigma}^2$

We have the distribution of $\hat{\alpha}$ and $\hat{\beta}$:

$$\begin{aligned}\hat{\beta} &\sim \mathcal{N}\left(\beta, \frac{\sigma^2}{S_{xx}}\right) \\ \hat{\alpha} &\sim \mathcal{N}\left(\alpha, \frac{\text{Var}(\hat{\beta}) \sum_{i=1}^n x_i^2}{n}\right).\end{aligned}$$

We write the following for the sample variance of α and β :

$$s_{\hat{\beta}}^2 = \frac{\hat{\sigma}^2}{S_{xx}}$$

$$s_{\hat{\alpha}}^2 = \frac{s_{\hat{\beta}}^2 \sum_{i=1}^n x_i^2}{n}.$$

In terms of t and χ^2 distributions:

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} = \frac{RSS}{\sigma^2} \sim \chi_{n-2}^2$$

$$\frac{\hat{\alpha} - \alpha}{s_{\hat{\alpha}}} \sim t_{n-2}$$

$$\frac{\hat{\beta} - \beta}{s_{\hat{\beta}}} \sim t_{n-2}.$$

In practice, we substitute α , β , and σ^2 for some hypothesised value. For example:

$$H_0 : \beta = 0 \Rightarrow \frac{\hat{\beta}}{s_{\hat{\beta}}} \sim t_{n-2}.$$

16 Confidence Intervals for α and β

16.1 Confidence Interval for α

The $\gamma\%$ confidence interval for α from a sample $\{x_1, x_2, \dots, x_n\}$ with response set $\{y_1, y_2, \dots, y_n\}$ is (c_L, c_U) where T_n^{-1} is the t_n quantile function (inverse CDF) and:

$$c_L = \hat{\alpha} - s_{\hat{\alpha}} T_{n-2}^{-1} \left(\frac{1+\gamma}{2} \right)$$

$$c_U = \hat{\alpha} + s_{\hat{\alpha}} T_{n-2}^{-1} \left(\frac{1+\gamma}{2} \right).$$

16.2 Confidence Interval for β

The $\gamma\%$ confidence interval for β from a sample $\{x_1, x_2, \dots, x_n\}$ with response set $\{y_1, y_2, \dots, y_n\}$ is (c_L, c_U) where T_n^{-1} is the t_n quantile function (inverse CDF) and:

$$c_L = \hat{\beta} - s_{\hat{\beta}} T_{n-2}^{-1} \left(\frac{1+\gamma}{2} \right)$$

$$c_U = \hat{\beta} + s_{\hat{\beta}} T_{n-2}^{-1} \left(\frac{1+\gamma}{2} \right).$$

17 Hypothesis Tests: Linear Regression

17.1 Hypothesis Test for β

17.1.1 Assumptions

For a sample $\{x_1, \dots, x_n\}$ randomly sampled with response set $\{y_1, \dots, y_n\}$ with Normal errors following a $\mathcal{N}(0, \sigma^2)$ and an expected value β_0 .

17.1.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \beta &= \beta_0 \\ H_1 : \beta &\neq \text{ or } > \text{ or } < \beta_0. \end{aligned}$$

17.1.3 Test Statistic

For $s_{\hat{\beta}}^2$ defined as follows:

$$s_{\hat{\beta}}^2 = \frac{\hat{\sigma}^2}{S_{xx}}, \quad (1)$$

we have the test statistic defined by:

$$T = \frac{(\hat{\beta} - \beta_0)}{s_{\hat{\beta}}} \sim t_{n-2}.$$

17.1.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let T_n^{-1} be the t_n quantile function (inverse CDF), γ be our significance level:

$$\begin{aligned} H_1 : \beta &\neq \beta_0 \Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = T_{n-2}^{-1}(1 - \frac{\gamma}{2})) \\ H_1 : \beta &> \beta_0 \Rightarrow C = [c^*, \infty) & (c^* = T_{n-2}^{-1}(1 - \gamma)) \\ H_1 : \beta &< \beta_0 \Rightarrow C = (-\infty, -c^*]. & (c^* = T_{n-2}^{-1}(1 - \gamma)) \end{aligned}$$

17.2 Hypothesis Test for α

17.2.1 Assumptions

For a sample $\{x_1, \dots, x_n\}$ randomly sampled with response set $\{y_1, \dots, y_n\}$ with Normal errors following a $\mathcal{N}(0, \sigma^2)$ and an expected value α_0 .

17.2.2 Hypotheses

We choose the hypotheses:

$$\begin{aligned} H_0 : \alpha &= \alpha_0 \\ H_1 : \alpha &\neq \text{ or } > \text{ or } < \alpha_0. \end{aligned}$$

17.2.3 Test Statistic

For $s_{\hat{\alpha}}^2$ defined as follows:

$$s_{\hat{\alpha}}^2 = s_{\hat{\beta}}^2 \left(\frac{\sum_{i=1}^n x_i^2}{n} \right), \quad (2)$$

we have the test statistic defined by:

$$T = \frac{(\hat{\alpha} - \alpha_0)}{s_{\hat{\alpha}}} \sim t_{n-2}.$$

17.2.4 Critical Region

For our observed test statistic t_{obs} , we want to generate a critical region of values. This region C changes based on our H_1 , let T_n^{-1} be the t_n quantile function (inverse CDF), γ be our significance level:

$$\begin{aligned} H_1 : \alpha \neq \alpha_0 &\Rightarrow C = (-\infty, -c^*] \cup [c^*, \infty) & (c^* = T_{n-2}^{-1}(1 - \frac{\gamma}{2})) \\ H_1 : \alpha > \alpha_0 &\Rightarrow C = [c^*, \infty) & (c^* = T_{n-2}^{-1}(1 - \gamma)) \\ H_1 : \alpha < \alpha_0 &\Rightarrow C = (-\infty, -c^*]. & (c^* = T_{n-2}^{-1}(1 - \gamma)) \end{aligned}$$