Random Samples

Definition 6.1

A set of random variables X_1, \ldots, X_n is called a Random Sample from a population if X_1, \ldots, X_n are mutually independent and each X_i has the same cdf F.

- * F describes the assumed distribution in the population.
- * Corresponding to F is a pdf (or pmf) f.
- * X_1, \ldots, X_n are Independent and Identically Distributed (*iid*).

Inference

* Joint pdf (pmf) of the sample

$$f(x_1,...,x_n \mid \boldsymbol{\theta}) = \prod_{i=1}^n f(x_i \mid \boldsymbol{\theta})$$

- * Usually the parameter vector θ is unknown.
- * Aim is to make inference about θ based on the observed sample x_1, \ldots, x_n .
- * Inference is based on statistics.

Statistics

Definition 6.2

Let X_1, \ldots, X_n be a random sample from an infinite population and let $T(x_1, \ldots, x_n)$ be a function mapping the support of X_1, \ldots, X_n , \mathcal{X}^n to \mathbb{R}^m where $m \leq n$. Then the random variable (or vector)

$$Y = T(X_1, \ldots, X_n)$$

is called a statistic. The distribution of the random variable Y is known as its sampling distribution.

* Since Y is a function of X_1, \ldots, X_n its sampling distribution can, in theory, be found from $f(x_1, \ldots, x_n \mid \theta)$.

Sample Mean and Variance

Definition 6.3

If X_1, \ldots, X_n is a random sample then the sample mean is

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

and the sample variance is

$$S^{2} = \frac{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}}{n-1}$$

The positive square root, S, of the sample variance is called the sample standard deviation.

* Observed values of X_1, \ldots, X_n are x_1, \ldots, x_n

* Observed values of \overline{X} and S^2 are

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
 $s^2 = \frac{\sum_{i=1}^{n} (x_i - \overline{x})^2}{n-1}.$

Linear Combinations

Theorem 6.1

Let X_1, \ldots, X_n be a sequence of random variables with finite means and variances and let a_1, \ldots, a_n be real constants. Then

$$E\left(\sum_{i=1}^{n} a_i X_i\right) = \sum_{i=1}^{n} a_i E(X_i)$$

$$Var\left(\sum_{i=1}^{n} a_i X_i\right) = \sum_{i=1}^{n} a_i^2 Var(X_i) + \sum_{j \neq i} a_i a_j Cov(X_i, X_j)$$

$$= \sum_{i=1}^{n} a_i^2 Var(X_i) + 2 \sum_{i>j} a_i a_j Cov(X_i, X_j)$$

Linear Combinations of Random Samples

Corollary 6.1.1

Let X_1, \ldots, X_n be a random sample from a distribution having finite mean, μ and finite variance, σ^2 , and let a_1, \ldots, a_n be real constants. Then

$$\mathsf{E}\left(\sum_{i=1}^{n} a_i X_i\right) = \mu \sum_{i=1}^{n} a_i$$
$$\mathsf{Var}\left(\sum_{i=1}^{n} a_i X_i\right) = \sigma^2 \sum_{i=1}^{n} a_i^2$$

Moments of Statistics

Lemma 6.1

Let X_1, \ldots, X_n be a random sample and let g be a function such that $Y = g(X_1)$ has finite mean and variance. Then

$$E\left(\sum_{i=1}^{n} g(X_i)\right) = n E\left(g(X_1)\right)$$
$$Var\left(\sum_{i=1}^{n} g(X_i)\right) = n Var\left(g(X_1)\right)$$

Convolutions

Theorem 6.2 (Bivariate Convolution)

If X and Y are independent random variables with pdfs f_X and f_Y respectively and Z = X + Y then the pdf of Z is

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(w) f_Y(z-w) dw.$$

Theorem 6.3 (General Convolution)

Let X_1, \ldots, X_n be a sequence of independent random variables such that X_i has pdf f_{X_i} and let $Z = \sum X_i$. Then the pdf of Z is

$$f_{Z}(z) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \left[f_{X_{1}}(w_{1}) \prod_{i=2}^{n-1} f_{X_{i}}(w_{i} - w_{i-1}) \right] \\ f_{X_{n}}(z - w_{n-1}) dw_{1} \cdots dw_{n-1}$$

Sample Mean

Theorem 6.4

Let X_1, \ldots, X_n be a random sample from a population with mean μ and finite variance σ^2 and let \overline{X} be the corresponding sample mean. Then

$$\mathsf{E}\left(\overline{X}\right) = \mu, \quad and \quad \mathsf{Var}\left(\overline{X}\right) = \frac{\sigma^2}{n}.$$

Theorem 6.5

Let X_1, \ldots, X_n be a random sample from a population with moment generating function $M_X(t)$ then the sampling distribution of the sample mean \overline{X} has moment generating function

$$M_{\overline{X}}(t) = \left[M_X\left(\frac{t}{n}\right)\right]^n$$

Estimators and Estimates

- * A statistic that is used to estimate a population quantity (parameter) θ is called an **estimator**.
- * The observed value of an estimator is called the estimate.

Definition 6.4

A statistic $T(X_1, ..., X_n)$ is said to be an **unbiased estimator** of the parameter θ if, and only if,

$$\mathsf{E}_{\theta}\left(T(X_1,\ldots,X_n)\right) = \theta$$

for all possible values of θ .

Theorem 6.6

Let X_1, \ldots, X_n be a random sample from a population with finite mean and variance μ and σ^2 . Then \overline{X} is an unbiased estimator of μ and S^2 is an unbiased estimator of σ^2 .

Samples from Exponential Family Distributions

Theorem 6.7

Suppose that X_1, \ldots, X_n is a random sample from a full (not curved) exponential family with common pdf (or pmf)

$$f(x \mid \boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp\left(\sum_{j=1}^{k} w_j(\boldsymbol{\theta})t_j(x)\right)$$

such that the set $\{w_1(\theta), \ldots, w_k(\theta)\}$ contains an open subset in $I\!\!R^k$

Define the statistics $T_j = \sum_{i=1}^n t_i(X_j)$ for j = 1, ..., k then the joint distribution of $T = (T_1, ..., T_k)$ is k-dimensional exponential family of the form

$$f_T(t_1, \ldots, t_k \mid \boldsymbol{\theta}) = h_1(t_1, \ldots, t_k) [c(\boldsymbol{\theta})]^n \exp\left(\sum_{i=1}^k w_i(\boldsymbol{\theta}) t_i\right)$$

Normal Random Samples

Theorem 6.8

Let X_1, \ldots, X_n be a sample from a $N(\mu, \sigma^2)$ population and let

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$$
 and $S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \overline{X})^2$

be the sample mean and variance. Then

(i) \overline{X} and S^2 are independent.

(ii) $\overline{X} \sim N(\mu, \sigma^2/n)$.

(iii) $(n-1)S^2/\sigma^2 \sim \chi^2_{n-1}$.

Pivotal Quantities

Definition 6.5

Suppose $X = (X_1, ..., X_n)$ is a sample from a population with cdf depending on some parameters θ . A quantity $R(X, \theta)$ which is a function of the data and the parameters is called a **pivotal quantity** (or simply a **pivot**) if the sampling distribution of R does not depend on the parameters θ .

* If
$$X_1, \ldots, X_n$$
 are *iid* N(μ, σ^2) then

$$Z = \frac{\sqrt{n}(\overline{X} - \mu)}{\sigma} \sim \mathsf{N}(0, 1)$$

* We shall see that another pivot in this situation is

$$T = \frac{\sqrt{n}(\overline{X} - \mu)}{S}$$

The Student's t Distribution

Theorem 6.9

If Z and X are two independent random variables with $Z \sim N(0,1)$ and $X \sim \chi^2_{\nu}$ then the random variable

$$T = \frac{Z}{\sqrt{X/\nu}}$$

has pdf given by

$$f_T(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\nu}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}} \quad \text{for } t \in \mathbb{R}.$$

The distribution with this pdf is called the Student's t distribution with ν degrees of freedom.

Properties of the t_{ν} Distribution

- * $E(T^r)$ exists if, and only if, $r < \nu$.
- * E(T) = 0 for $\nu > 1$.
- * $Var(T) = \frac{\nu}{\nu 2}$ for $\nu > 2$.
- * Suppose that $X \sim t_1$ then

$$f_X(x) = \frac{1}{\pi(1+x^2)} \quad -\infty < x < \infty.$$

This is called the standard Cauchy distribution.

Connection between T and Normal Distributions

Theorem 6.10

Suppose that $T_1, T_2, ...$ is a sequence of random variables such that $T_{\nu} \sim t_{\nu}$ and $Z \sim N(0, 1)$. Then

 $\mathsf{P}(T_{\nu} \leq x) \rightarrow \mathsf{P}(Z \leq x) \text{ as } \nu \rightarrow \infty \text{ for any } x \in \mathbb{R}$

Theorem 6.11

Suppose that X_1, \ldots, X_n is a random sample from a Normal (μ, σ^2) population and that \overline{X} and S^2 are the sample mean and sample variance. Then

$$T = \frac{\sqrt{n}(\overline{X} - \mu)}{S} \sim t_{n-1}$$

Two-Sample Inference

Theorem 6.12

Suppose that X_1, \ldots, X_n and Y_1, \ldots, Y_m are independent random samples from normal populations with parameters (μ_X, σ^2) and (μ_Y, σ^2) respectively. Let \overline{X} and \overline{Y} be the sample means and S_X^2 and S_Y^2 be the sample variances. Define the **pooled variance** estimate

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

Then

(i)
$$\frac{(n+m-2)S_p^2}{\sigma^2} \sim \chi_{n+m-2}^2$$

(ii)
$$T = \frac{(\overline{X} - \overline{Y}) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

Snedecor's *F* **Distribution**

Definition 6.6

A random variable Y is said to have and F distribution with p numerator degrees of freedom and q denominator degrees of freedom if, and only if, its pdf is given by

$$f_Y(y) = \frac{\Gamma\left(\frac{p+q}{2}\right)}{\Gamma\left(\frac{p}{2}\right)\Gamma\left(\frac{q}{2}\right)} \left(\frac{p}{q}\right)^{p/2} \frac{y^{(p/2)-1}}{[1+(p/q)y]^{(p+q)/2}} \text{ for } 0 < y < \infty.$$

Theorem 6.13

Suppose that X_1 has a Chi-squared(p) distribution, X_2 has a Chi-squared(q) distribution and X_1 and X_2 are independent. Then the random variable

$$Y = \frac{X_1/p}{X_2/q}$$

has an F distribution with p and q degrees of freedom.

Comparison of 2 Normal Variances

Theorem 6.14

Suppose that X_1, \ldots, X_n and Y_1, \ldots, Y_m are independent random samples from normal populations with parameters (μ_X, σ_X^2) and (μ_Y, σ_Y^2) respectively. Let S_X^2 and S_Y^2 be the sample variances. Then

$$\frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2} \sim F_{n-1,m-1}$$

Order Statistics

Definition 6.7

Let X_1, \ldots, X_n be a random sample then the order statistics of the sample are denoted $X_{(r)}$, $r = 1, \ldots, n$ where

$$X_{(1)} \leqslant X_{(2)} \leqslant \cdots \leqslant X_{(n)}.$$

Theorem 6.15

Let X_1, \ldots, X_n be a random sample from a distribution with cdf F_X . Then the cdf of the sample maximum, $X_{(n)}$, is

$$F_{X_{(n)}}(x) = [F_X(x)]^n$$
.

and that for the minimum, $X_{(1)}$, is

$$F_{X_{(1)}}(x) = 1 - [1 - F_X(x)]^n$$
.

Distribution of Order Statistics (Discrete)

Theorem 6.16

Let X_1, \ldots, X_n be a random sample from a discrete distribution on the values $x_1 < x_2 < \cdots$. Let the common probability mass function of the random variables be $P(X = x_i) = p_i$ with corresponding cdf

$$\mathsf{P}(X \leqslant x_i) = P_i = \sum_{k=1}^i p_k$$

and let us define $P_0 = 0$.

If $X_{(r)}$ is the rthorder statistic of the sample then

$$\mathsf{P}(X_{(r)} \leq x_i) = \sum_{k=r}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k}$$

and
$$P(X_{(r)} = x_i) = \sum_{k=r}^{n} {n \choose k} \left[P_i^k (1 - P_i)^{n-k} - P_{i-1}^k (1 - P_{i-1})^{n-k} \right]$$

Distribution of Order Statistics (Continuous)

Theorem 6.17

Let X_1, \ldots, X_n be a random sample from a continuous distribution with pdf f_X and cdf $F_X(x)$ and let $X_{(r)}$ be the r^{th} order statistic. Then the pdf of $X_{(r)}$ is

$$f_{X_{(r)}}(x) = \frac{n!}{(r-1)!(n-r)!} f_X(x) \left[F_X(x)\right]^{r-1} \left[1 - F_X(x)\right]^{n-r}$$

٠

Joint Distribution of Two Order Statistics (Continuous)

Theorem 6.18

Let X_1, \ldots, X_n be a random sample from a continuous distribution with pdf f_X and cdf $F_X(x)$ and let $X_{(r)}$ and $X_{(s)}$ be two order statistics with r < s. Then the joint pdf of $X_{(r)}$ and $X_{(s)}$ is

$$f_{X_{(r)},X_{(s)}}(u,v) = \frac{n!}{(r-1)!(s-r-1)!(n-s)!} f_X(u) f_X(v)$$
$$\times [F_X(u)]^{r-1} [F_X(v) - F_X(u)]^{s-r-1} [1 - F_X(v)]^{n-s}$$

for $-\infty < u < v < \infty$.

Joint Distribution of All Order Statistics (Continuous)

Theorem 6.19

Let X_1, \ldots, X_n be a random sample from a continuous distribution with pdf f_X and let $X_{(1)}, \ldots X_{(n)}$ be the order statistics. Then the joint pdf of all of the order statistics is

$$f_{X_{(1)},\dots,X_{(n)}}(x_1,\dots,x_n) = n! f_X(x_1)\cdots f_X(x_n)$$

for $-\infty < x_1 < \cdots < x_n < \infty$.