

# Métodos Estadísticos

Salvador Pintos  
Instituto de Cálculo Aplicado

julio 2000

# Índice General

1	Revisión de conceptos básicos	3
1.1	Variables aleatorias	3
1.1.1	Estadística descriptiva en el SAS	3
1.1.2	Variables Aleatorias, Valor Esperado y Varianza	4
1.1.3	Valores esperados y varianzas	5
1.2	Distribuciones discretas	5
1.2.1	Distribución binomial	5
1.2.2	Distribución Hipergeométrica	6
1.2.3	Distribución binomial negativa	6
1.2.4	Distribución de Poisson	6
1.3	Distribuciones continuas	7
1.3.1	La Distribución Normal	7
1.3.2	Distribución Gamma	8
1.3.3	Distribución Exponencial	8
1.3.4	Distribución Beta	8
1.3.5	La Distribución Uniforme	9
1.4	Simulación de variables aleatorias con el SAS	10
2	Inferencia	11
2.1	Estimación puntual	11
2.1.1	Muestreo aleatorio simple	11
2.1.2	Propiedades de los estimadores puntuales	12
2.1.3	Ley de los Grandes Números y el Teorema Central del Límite	13
2.1.4	Distribuciones asociadas a la Normal en el muestreo	15
2.2	Estimación por intervalos	17
2.2.1	Intervalos de la media y de la varianza	18
2.2.2	Tamaño muestral	19
2.2.3	Intervalo de confianza para proporciones	20

2.3	Prueba de hipótesis . . . . .	21
2.3.1	Conceptos básicos . . . . .	22
2.3.2	Como programar las pruebas . . . . .	24
2.3.3	Pruebas para dos poblaciones . . . . .	25
2.3.4	Pruebas de ajuste de distribuciones . . . . .	29
2.3.5	Prueba de independencia . . . . .	30

# Capítulo 1

## Revisión de conceptos básicos

### 1.1 Variables aleatorias

#### 1.1.1 Estadística descriptiva en el SASS

Una descripción informativa del conjunto de datos de una muestra de una variable está dada por el histograma de frecuencias relativas. La habilidad para detectar patrones o tipos de distribución a partir de un histograma depende de la elección adecuada de las clases. El PROC CHART es el procedimiento idóneo para hacer estos histogramas. Se sugiere hacer uso de la opción midpoints para controlar el número de clases del histograma.

Existen dos medidas de interés para cualquier conjunto de datos  $\{x_1, \dots, x_n\}$ : la localización de su centro y la variabilidad.

Hay, principalmente, tres medidas de tendencia central:

La media muestral:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

La mediana muestral: si se ordena el conjunto de datos la mediana divide la muestra en dos.

La moda muestral: es el valor más frecuente de la muestra.

Es la media, por razones teóricas, la que se usa fundamentalmente.

En cuanto a la variabilidad, la varianza muestral es la medida estadísticamente más útil:

$$S^2 = \frac{\sum_{i=1}^n (x_i - m)^2}{(n - 1)}$$

A menudo se prefiere la desviación standard, que es su raíz cuadrada, porque se expresa en las mismas unidades físicas que la media y las observaciones. El PROC UNIVARIATE del SAS permite calcular esas medidas y otras: cuantiles, observaciones máximas y mínimas, etc.

### 1.1.2 Variables Aleatorias, Valor Esperado y Varianza.

Cuando se consideran variables económicas o sociales es necesario admitir que son esencialmente Variables Aleatorias y que en consecuencia tienen asociada una estructura de probabilidad que se caracteriza por la distribución de probabilidad.

Una Variable Aleatoria  $X$  es Discreta si el conjunto de valores que toma es...nito; si es infinito puede ordenarse en una secuencia que se corresponda con los naturales.

Su distribución de probabilidad está dada por:  $P\{X = x_i\} = p_i$ .

Una Variable Aleatoria  $X$  es Continua si el conjunto de valores que toma es uno o más intervalos de la recta real.

Su distribución de probabilidad está caracterizada por la función de densidad  $f$ .

Definición 1  $f$  es la función de densidad de la variable aleatoria  $X$  si  $f(x) \geq 0$  y para cada intervalo  $[a, b]$ :

$$P\{a \leq X \leq b\} = \int_a^b f(x) dx$$

La Función de Distribución Acumulada es  $F(x) = P\{X \leq x\}$ :

Con el propósito de calcular probabilidades, el SAS tabula, para las distintas distribuciones, esta Función de Distribución Acumulada.

Si  $X$  es discreta y toma valores ordenados  $0, 1, 2, \dots$  etc. entonces:  $P\{X = k\} = F(k) - F(k - 1)$

Para una Variable Aleatoria  $X$  es necesario establecer su valor medio (Valor Esperado), y cómo se dispersa respecto de su valor medio (Varianza).

Definición 2 Si  $X$  es continua Valor Esperado  $E(X)$  es:

$$E(X) = \int_{-\infty}^{\infty} xf(x) dx$$

la Varianza  $V(X)$ :  $V(X) = E(X - E(X))^2$  y la Desviación Standard  $\sigma = \sqrt{V(X)}$

### 1.1.3 Valores esperados y varianzas

En lo que sigue usaremos  $E$ ,  $V$  y  $\sigma$  como símbolos de Valor Esperado, Varianza y Desviación Estándar.

Proposición 3 Si  $X$  e  $Y$  son variable aleatorias y  $k$  es una constante entonces:

$$E(kX) = kE(X) \quad V(kX) = k^2V(X) \quad \sigma(kX) = |k|\sigma(X)$$

Además si  $X$  e  $Y$  son independientes:  $V(X + Y) = V(X) + V(Y)$

$$\sigma(X + Y) = \sqrt{V(X) + V(Y)}$$

Ejemplo Si el peso de un níspero se distribuye normal con valor esperado

120 grs y desviación estándar 15 grs, la balsa también normal con peso 5 grs y desviación estándar 1 gr y el precio es de Bs 100 por kilo. Hallar las siguientes probabilidades:

- Q ue una balsa con 3 nísperos pese más de 400 grs.
- Q ue 10 nísperos sin balsa pesen más de 1300 grs.
- Q ue 1 níspero cueste más de 13 Bs.

## 1.2 Distribuciones discretas

### 1.2.1 Distribución binomial

Características:

1. Hay  $n$  ensayos independientes.
2. El resultado del ensayo es éxito (E) o fracaso (F).
3. La probabilidad  $p$  de éxito es constante en los ensayos.
4. Distribución asociada al muestreo con reemplazo

Parámetros:  $n$ ,  $p$ .

Variable aleatoria: número  $x$  de éxitos en los  $n$  ensayos.

Función de distribución acumulada en SAS: `PROBBNML(p, n, x)`

### 1.2.2 Distribución Hipergeométrica

Características:

1. D es una población de tamaño  $N$ ,  $P$ ,  $K$  elementos son del tipo  $A$ . Si se selecciona una muestra aleatoria de  $n$  elementos ¿cuál es la probabilidad de que en ésta se hallen  $x$  elementos del tipo  $A$ ?

2. Muestreo sin reposición

Parámetros:  $N$ ,  $P$ ,  $K$ ,  $n$ .

Variable aleatoria: número  $x$  de elementos  $A$  en la muestra

Obsérvese que  $0 \leq x \leq K$  y  $0 \leq n - x \leq N - P$ . Luego los posibles valores de la variable aleatoria  $x$  están restringidos al intervalo  $\max(0; n + K - N) \leq x \leq \min(n, K)$

Sea  $p = K/N$  la proporción de elementos del tipo  $A$  en la población:

Valor esperado y Varianza  $E(X) = np$   $V(X) = np(1-p)[(N-1)/(N+1)]$

Función de distribución acumulada en SAS: `PROBHYPR(N, P, K, n, x)`

Cuando la fracción de muestreo  $n/N$  es pequeña, la hipergeométrica se puede aproximar por la binomial con parámetros  $p$  y  $n$ .

### 1.2.3 Distribución binomial negativa

Características:

1. El resultado del ensayo es éxito ( $E$ ) o fracaso ( $F$ ).
2. La probabilidad  $p$  de éxito es constante en los ensayos.
3. Se realizan ensayos independientes consecutivos hasta obtener  $k$  éxitos.

Parámetros:  $k$ ,  $p$ .

Variable aleatoria: número  $x$  de fracasos antes del  $k$  éxito

Función de distribución acumulada en SAS: `PROBNGB(p, k, x)`

La distribución geométrica es un caso particular de ésta cuando  $k=1$ , es decir la variable aleatoria de la distribución geométrica representa el número de fracasos antes de obtener el primer éxito

### 1.2.4 Distribución de Poisson

Características:

1. Eventos que ocurren con velocidad constante en el tiempo (número de carros que llegan a un autódromo), o en el espacio (número de fallas por metro en un rollo de tela, etc.).

2. Se utiliza como aproximación al modelo binomial cuando  $n$  es grande y  $p$  pequeño (ley de los sucesos raros).  $p < 0.1$   $np < 5$ , tomando como parámetro  $\lambda = np$ .

Parámetros:  $\lambda$  promedio de ocurrencias del suceso por unidad de tiempo o espacio

Función de distribución acumulada en SAS: `POISSON( $\lambda$ , x)`

## 1.3 Distribuciones continuas

### 1.3.1 La Distribución Normal

La distribución Normal es indudablemente la distribución continua fundamental, tanto por sus aplicaciones como por el rol que juega dentro de la Teoría Estadística. Es la piedra angular de la Inferencia ya que muchas estadísticas muestrales tienden hacia la distribución Normal cuando el tamaño de la muestra crece.

Se afirma que una variable aleatoria  $X$  es Normal  $N(\mu, \sigma^2)$  si su función de densidad está dada por:

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2} \left(\frac{x - \mu}{\sigma}\right)^2\right]$$

$-\infty < x < \infty \quad -\infty < \mu < \infty \quad \sigma > 0$

Los parámetros  $\mu$  y  $\sigma$  son, además, el Valor Esperado y la Desviación Standard de  $X$ . Si  $X$  es  $N(\mu, \sigma^2)$  entonces,  $Z = (X - \mu)/\sigma$ , es  $N(0; 1)$  y se le designa como Normal Standard. Esta propiedad permite relacionar la función de distribución acumulada de  $X$  con la de  $Z$ , ya que luego es suficiente con tabular las probabilidades de la distribución Normal Standard. La distribución acumulada de  $Z$  en el SAS es `PROBNORM(z)` y su inversa `PROBIT(p)`.

Ejemplo `prb(z = 2) = PROBNORM(2)`

Ejemplo ¿Cuál es el  $z$  tal que la `Prb(z < z) = 0.8`? `z = PROBIT(0.8)`.

Ejemplo Si el salario de un obrero de la Industria Petrolera se distribuye Normal con  $\mu =$  Bs 17.000 y  $\sigma =$  Bs 3.500, ¿cuál es la probabilidad de que un obrero gane más de Bs 20.500? `p = PROBNORM((20500-17000)/3500)` es la probabilidad de que gane menos de 20500 y `1-p`, su complemento, es la probabilidad de que gane más de 20500.

Ejercicio Hacer un programa en SAS para calcular esta probabilidad.

La distribución Normal tiene además la propiedad fundamental siguiente



Si una variable aleatoria es la combinación lineal de variables aleatorias normales independientes, entonces es también una variable aleatoria normal.

### 1.3.2 Distribución Gamma

La distribución Gamma es una distribución continua con dos parámetros:

$\mu$  = Parámetro de escala,  $\alpha$  = parámetro de forma

Su función de densidad es:

$$f(x; \mu; \alpha) = \frac{x^{\alpha-1} e^{-\frac{x}{\mu}}}{\Gamma(\alpha) \mu^\alpha}$$

La forma depende de  $\alpha$ . Si  $\alpha > 1$  la densidad crece y luego decrece con la moda en  $\frac{\alpha-1}{\mu}$ .

Cuando  $\alpha > 1$  a la distribución se le llama exponencial. Esta distribución juega un rol esencial en la teoría de colas ya que representa el tiempo entre dos llegadas consecutivas, cuando las llegadas son independientes.

La Distribución Gamma en el SAS.

El SAS tabula solamente la distribución Gamma estándar ( $\mu = 1$ ). La función de distribución acumulada es:  $\text{PROBGAM}(x, \alpha)$  y su inversa:  $\text{GAMINV}(\text{prcb}, \alpha)$ .

¿Cómo obtener entonces la función de distribución acumulada de una Gamma de parámetros  $\alpha, \mu$ ? Se logra a partir de una Gamma standard con  $\mu = 1$ . Si se define como  $\text{PROBGAM2}(x, \alpha, \mu)$  a la función con los dos parámetros se cumple que:

$$\begin{aligned} \text{PROBGAM2}(x, \alpha, \mu) &= \text{PROBGAM}\left(\frac{x}{\mu}, \alpha\right) \\ \text{GAMINV2}(\text{prcb}, \alpha, \mu) &= \mu * \text{GAMINV}(\text{prcb}, \alpha) \end{aligned}$$

### 1.3.3 Distribución Exponencial

La distribución Exponencial de parámetro  $\mu$  es una Gamma generalizada con  $\alpha = 1$  entonces:

$$\begin{aligned} \text{PROBEXP}(x; \mu) &= \text{PROBGAM2}(x; 1; \mu) = \text{PROBGAM}\left(\frac{x}{\mu}; 1\right) \\ \text{EXPINV}(\text{prcb}; \mu) &= \mu * \text{GAMINV}(\text{prcb}; 1) \end{aligned}$$

### 1.3.4 Distribución Beta

La distribución Beta es una distribución continua con dos parámetros:  $\alpha$  y  $\beta$ , ambos parámetros de forma.

Se la utiliza para representar variables aleatorias cuyos valores se encuentran

restringidos a intervalos de longitud infinita. Ejemplos: PERT, evaluación de programas, límites estadísticos de tolerancia, estadística del orden, etc.

Su función de densidad es 0 en todos lados salvo en el intervalo  $[0,1]$  donde está definida por:

$$f(\alpha; \beta; x) = \frac{\Gamma(\alpha + \beta) x^{\alpha-1} (1-x)^{\beta-1}}{\Gamma(\alpha) \Gamma(\beta)} \quad 0 < x < 1, \quad \alpha > 0, \quad \beta > 0$$

La forma depende de la posición de los parámetros  $\alpha$  y  $\beta$  respecto de 1.

Ejercicio: ejecute el programa betapl y observe la forma de la curva en función de los parámetros.

Si  $\alpha$  y  $\beta$  son mayores que 1 la moda es:  $(\alpha - 1) = (\beta - 1)$

La distribución Beta en el SAS

PROBBETA( $x, \alpha, \beta$ ) es la función de distribución acumulada definida para todo  $x$  en  $(0;1)$ ; BETAINV( $p, \alpha, \beta$ ) su función inversa.

### 1.3.5 La Distribución Uniforme

La distribución uniforme o rectangular tiene densidad constante en un intervalo  $[a,b]$  y vale 0 fuera del mismo. En consecuencia está dada por:

$$f(x) = \frac{1}{b-a} \text{ en } [a,b]$$

$$E(X) = \frac{a+b}{2} \quad \text{VAR}(X) = \frac{(b-a)^2}{12}$$

Esta distribución se aplica en teoría de errores y en la generación de observaciones simuladas que sigan una distribución dada.

### Aproximación Normal de la Binomial

Cuando en el modelo Binomial  $n$  es grande, y tanto  $np$  como  $n(1-p)$  son mayores que 5, se puede aproximar la binomial con una Normal de igual esperanza y varianza. Esto es con una  $N[np; np(1-p)]$ .

Sin embargo se tendrá en cuenta la siguiente:

$\text{Pr}(X = n)$  se aproxima por  $\text{Pr}(n-0.5 < X < n+0.5)$

$\text{Pr}(n_1 \leq X \leq n_2)$  se aproxima por:

$\text{Pr}(n_1 - 0.5 \leq X \leq n_2 + 0.5)$

Ejemplo: Si una coneja tiene 50 hijos ¿cuál es la probabilidad de tener exactamente 25 hembras? ¿Cuál la de tener entre 20 y 30 machos?

Si  $p$  es pequeño ( $p < 0.1$ ), y  $np < 5$ , la binomial se aproxima por la distribución de Poisson con valor  $\lambda = np$ .

## 1.4 Simulación de variables aleatorias con el SAS

El SAS dispone de varias funciones para generar variables aleatorias que sigan una distribución dada. Dos muy útiles son:

$Z = \text{RANNOR}(\text{semilla})$ , genera valores de una Normal STANDARD,  $N(0,1)$

$U = \text{RANUNI}(\text{semilla})$ , genera valores de una Uniforme en el intervalo  $[0,1]$

Puesto que en esencia son pseudo aleatorias, todas dependen de una semilla inicial que genera la secuencia aleatoria. Si se cambia la semilla se genera otra secuencia de números independientes de la secuencia anterior. Elija para la semilla cualquier entero positivo para garantizar aleatoriedad, si elige 0 el SAS la genera a partir del tiempo del computador. Por ejemplo son válidas las siguientes semillas: SE= 3719 241; SE1= 258 583.

tabla de generación aleatoria con el sas

Distribución	parámetros	generación aleatoria
Uniforme	a, b	$u = a + (b-a) * \text{ranuni}(se)$
Normal	$\mu$ , $\sigma$	$nor = \mu + \sigma * \text{rannor}(se)$
Gamma	$\theta$ , $\alpha$	$gam = \theta * \text{rangam}(se, \alpha)$
Poisson	$\lambda$	$poisson = \text{ranpoi}(se, \lambda)$
Binomial	n, p	$binomi = \text{ranbin}(se, n, p)$
Geométrica	p	$ge = \lceil \text{cor}(j, \text{ranexp}(se) / \log(1-p)) \rceil$

Para una variable aleatoria continua  $X$  cuya Función de distribución  $F(x)$  es conocida, y donde su inversa es  $INV F(u)$ , es posible generar valores aleatorios de  $X$  mediante el siguiente procedimiento:

1. genere un valor distribuido uniformemente en  $[0,1]$
2. halle  $X = INV F(u)$

Los valores de  $RANUNI$ , básicos en la generación pseudo aleatoria, son generados por el método de congruencias donde el módulo es  $2^{31} - 1$  y el multiplicador 39 720 409 4. La semilla debe ser menor que el módulo.

## Capítulo 2

# Inferencia

### 2.1 Estimación puntual

#### 2.1.1 Muestreo aleatorio simple

**Definición 4** Dada una Variable aleatoria  $X$ , un conjunto  $\{X_1; X_2; \dots; X_n\}$  de variables es una Muestra Aleatoria Simple (M.A.S), de  $X$  si:

1. Cada  $X_i$  tiene la misma distribución que la  $X$ .
2. Las variables  $X_1; X_2; \dots; X_n$  son estadísticamente independientes entre sí.

En consecuencia, una muestra aleatoria simple (M.A.S), es un caso particular de distribución multivariada donde las variables  $\{X_1; X_2; \dots; X_n\}$  son independientes y todas tienen por distribución común la distribución de  $X$ . A menudo se confunde este concepto ya que también se le llama muestra al resultado de observar el valor que toman las  $n$  variables  $\{X_1; X_2; \dots; X_n\}$ . Sin embargo, y para que no queden dudas, las propiedades estadísticas que aquí se estudian se refieren a la M.A.S como variable aleatoria multivariada y no como al conjunto de números resultantes de la observación de la misma.

**Definición 5** Un estadístico  $T$  es una función de la muestra  $\{X_1; X_2; \dots; X_n\}$ . Es una nueva variable aleatoria uni-dimensional construida a partir de una muestra aleatoria simple.

Todo estadístico, por ser una variable aleatoria de una dimensión, tiene su distribución de probabilidad y suele ser más sencillo operar con él que con la muestra como tal. Por ejemplo, el estadístico más importante es la media muestral:  $m = \frac{\sum_{i=1}^n x_i}{n}$  y cuando se hace referencia a las propiedades de la media muestral se sobreentiende que es sobre el estadístico como variable

aleatoria

Otros estadísticos son:  $\max\{X_1; X_2; \dots; X_n\}$ ;  $\min\{X_1; X_2; \dots; X_n\}$ ; rango =  $\max_i - \min_i$ :

El propósito del muestreo estadístico es hacer inferencias acerca de la población. Más específicamente: estimar los parámetros desconocidos de la variable aleatoria  $X$  que representa a la población.

Una vez identificados los parámetros, se determina totalmente la variable aleatoria en el sentido de que se determina totalmente su distribución de probabilidad. Para ello se usan estadísticos adecuados, útiles, cuyas distribuciones se conocen, como la media muestral, la varianza muestral, el máximo de la muestra, etc.

**Definición 6** Si se emplea un estadístico  $T$  para estimar un parámetro  $\mu$ , entonces se dice que  $T$  es un estimador de  $\mu$ , y al resultado de  $T$ , al calcularlo sobre los valores específicos de la muestra, se le llama estimación de  $\mu$ .

### 2.1.2 Propiedades de los estimadores puntuales

Una estimación obtenida a partir del valor de un estimador, se la define como estimación puntual, por oposición a la estimación por intervalos que se verá más adelante. Algunas de las propiedades deseables más importantes de los estimadores puntuales se incluyen a continuación.

**Definición 7**  $T$  es un estimador insesgado de  $\mu$  si  $E(T) = \mu$  para todos los posibles valores de  $\mu$ . Se entiende por sesgo del estimador a la diferencia  $E(T) - \mu$ . Si  $T$  es insesgado su distribución se encuentra centrada en torno a  $\mu$ .

Es razonable esperar que un buen estimador  $T_n$ , (donde se ha agregado la  $n$  en la notación para hacer énfasis en la dependencia del tamaño muestral  $n$ ), se concentre en torno al parámetro  $\mu$  a medida que la muestra crece (o sea cuando la información aumenta).

**Definición 8**  $T_n$  es un estimador consistente de  $\mu$  si:

$$\lim_{n \rightarrow \infty} P(\text{rdo}(jT_n - \mu) \leq \epsilon) = 1 \quad (2.1)$$

para todos los valores de  $\mu$  y donde  $\epsilon$  es un número positivo arbitrario

El requisito establecido por el límite anterior constituye lo que se denomina convergencia en probabilidad. Luego  $T_n$  es un estimador consistente de  $\mu$  si converge en probabilidad a  $\mu$ . La consistencia significa la concentración del estimador en torno  $\mu$  a medida que la muestra crece.

Es esta propiedad de consistencia de un estimador lo que permite conformarse con el valor observado de un estimador y asumirlo como representativo del verdadero valor del parámetro  $\mu$ . Y a que, el error que se comete al sustituir el verdadero valor del parámetro por su estimación tiene una alta probabilidad de ser menor que  $\epsilon$ .

La varianza de un estimador insesgado es la cantidad más importante para decidir qué tan bueno es el estimador.

**Definición 9** Dados dos estimadores insesgados  $T_a$  y  $T_b$  de  $\mu$ , se dice que  $T_a$  es más eficiente que  $T_b$  si:

$$\text{Var}(T_a) < \text{Var}(T_b) \quad (2.2)$$

**Definición 10**  $T$  es un estimador insesgado de varianza mínima de  $\mu$  si es más eficiente que cualquier otro estimador insesgado para todos los valores posibles de  $\mu$ .

Es decir, si es el que tiene mínima varianza entre todos los estimadores insesgados de  $\mu$ .

### 2.1.3 Ley de los Grandes Números y el Teorema Central del Límite

La media muestral es el estadístico más usado en los procesos de inferencia:

$$\bar{m} = \frac{\sum_{i=1}^n X_i}{n} \quad (2.3)$$

Este estadístico tiene propiedades resaltantes que justifican su relevancia en el muestreo. Estas se resumen en dos grandes propiedades:

$\supset$  Ley de los Grandes Números

$\supset$  Teorema Central del Límite

Ambas válidas cuando la variable aleatoria  $X$  tiene valor esperado y varianza finitos.

Considérese una M.A.S. de tamaño  $n$ , de una variable aleatoria  $X$  de valor esperado  $\mu$  y desviación estándar  $\sigma$ , entonces por las propiedades ya establecidas sobre la suma de variables aleatorias independientes de la ecuación 2.3 se deduce:

$$E(\bar{m}) = \mu \quad \text{Var}(\bar{m}) = \frac{\sigma^2}{n} \quad (2.4)$$

Luego  $\bar{m}$  es un estimador insesgado de  $\mu$ . La media muestral tiene el mismo valor esperado  $\mu$  pero está más concentrada en torno a  $\mu$ , ya que la desviación estándar disminuye con la raíz cuadrada de  $n$  y  $\frac{\sigma^2}{n} \rightarrow 0$  cuando  $n \rightarrow \infty$ . Esta concentración se refleja en la siguiente propiedad que es el fundamento teórico para la estimación de la media poblacional por la media muestral:

**Teorema 11 Ley de los Grandes Números:**  
La media muestral  $\bar{m}$  es un estimador consistente de  $\mu$ .

La ley de los Grandes Números indica que  $\bar{m}$  se concentra entorno a  $\mu$ , pero no indica cómo es la distribución de la media muestral. El siguiente teorema tiene ese propósito:

**Teorema 12 Teorema central del límite**  
Cuando  $n \rightarrow \infty$  la distribución de la media muestral,  $\bar{m}$ , se aproxima a la de una normal; más precisamente el estadístico de la media estandarizada:

$$Z = \frac{\bar{m} - \mu}{\frac{\sigma}{\sqrt{n}}} \quad (2.5)$$

tiene una distribución que se aproxima a la normal estándar cuando  $n \rightarrow \infty$ .

Este último teorema es el centro de la estadística y el que le otorga a la normal su rol preponderante; ya que para muestras aleatorias simples de tamaño grande, de cualquier variable aleatoria de distribución no especificada, se puede trabajar con la hipótesis de que  $\bar{m}$  es normal, independientemente de la distribución de origen.

**Proposición 13 La varianza muestral**

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{m})^2}{(n-1)} \quad (2.6)$$

es un estimador insesgado de la varianza de la población, ( $E(s^2) = \sigma^2$ ), y además, consistente.

Pero no existe para ésta un teorema de convergencia como el teorema central del límite, tal como existe para la media muestral.

#### 2.1.4 Distribuciones asociadas a la Normal en el muestreo

En lo que sigue se considerará una muestra aleatoria simple  $(X_1; X_2; \dots; X_n)$  de una variable aleatoria Normal  $X \sim N(\mu; \sigma^2)$ . Luego, su media muestral  $\bar{X}$  se distribuye también Normal pero  $N(\mu; \sigma^2/n)$  cualquiera sea el tamaño de la muestra.

La distribución Chi-Cuadrado ( $\chi_p^2$ )

$$\chi_p^2 = \sum_{k=1}^p Z_k^2 \quad (2.7)$$

donde  $Z_1; \dots; Z_p$  es una muestra aleatoria simple de Normales estándar ( $N(0; 1)$ ).

Proposición 14 La distribución de la siguiente función de la varianza muestral  $S^2$ :

$$\frac{(n-1)S^2}{\sigma^2} \quad (2.8)$$

es una  $\chi^2$  (Chi-Cuadrado) con  $n-1$  grados de libertad.

Esta distribución que depende de un solo parámetro  $g$  (grados de libertad), es un caso particular de la distribución GAMMA donde el parámetro de forma  $\alpha = g/2$ , y el parámetro de escala  $\mu = 2$ :

Su distribución acumulada en  $S^2$  es  $PROBCH I(x, g)$ , y su inversa  $CINV(p, g)$ , que es la más útil en el muestreo, donde  $g$  significa grados de libertad.

Si  $X$  es una CHI-SQUARE entonces:  $E(X) = g$  y  $V(X) = 2g$

De la ecuación 2.8 se deduce que  $S^2 = \frac{\sigma^2}{n-1} \chi^2(n-1)$ , luego  $E(S^2) = \sigma^2$  y  $V(S^2) = 2\sigma^4/(n-1)$

Se prueba así que  $S^2$  es un estimador insesgado y consistente de la varianzaplacional  $\sigma^2$ , ya que la varianza tiende a 0 cuando el tamaño muestral crece.



## La distribución t de Student

Definición 15 La distribución  $t_p$  es el cociente entre las variables aleatorias independientes siguientes:

$$\frac{N(0;1)}{\sqrt{\frac{\chi_p^2}{p}}} \quad (2.9)$$

Esta distribución depende solamente de un parámetro  $p$ , (grados de libertad).  $E(t_p) = 0$  y  $Var(t_p) = \frac{p}{p-2}$  si  $p > 2$ . Cuando los grados de libertad crecen ( $p \rightarrow \infty$ ) la distribución  $t_p$  tiende a la normal estándar.

## La distribución F de Fischer

Definición 16 Si  $X$  e  $Y$  son variables aleatorias con distribuciones Chi-cuadrado de  $n_x$  y  $n_y$  grados respectivamente, se dice que la variable

$$F = \frac{\frac{X}{n_x}}{\frac{Y}{n_y}} \quad (2.10)$$

tiene una distribución F con  $n_x$  y  $n_y$  grados de libertad.

Esta distribución depende de dos parámetros, ( $g_n$  y  $g_d$ ), grados de libertad del numerador y denominador.

Su distribución acumulada en SAS es  $PROBF(x, g_n, g_d)$  y su inversa (que es la más útil en el muestreo)  $FINV(p, g_n, g_d)$ . El archivo `dist_ f.sas` permite graficar la distribución para distintos valores de los parámetros. Obsérvese la asimetría positiva para cualquier valor de los parámetros.

Su importancia reside en que

Proposición 17 si  $\{x_1, \dots, x_{n_x}\}$  es una muestra de una variable  $X \sim (1, \frac{1}{n_x})$  y  $\{y_1, \dots, y_{n_y}\}$  es una muestra de una variable  $Y \sim (1, \frac{1}{n_y})$ , entonces:

$$F_m = \frac{\frac{s_x^2}{\sigma_x^2}}{\frac{s_y^2}{\sigma_y^2}} \quad (2.11)$$

tiene por distribución una F con  $n_x - 1$  y  $n_y - 1$  grados de libertad

Además su valor esperado está dado por:  $E(F_{n_x-1, n_y-1}) = \frac{n_y}{n_y-2}$   $n_y > 2$

Si las variables  $X$  e  $Y$  tienen igual varianza (aunque sea desconocida), el valor muestral  $F_m$  depende sólo de las varianzas muestrales, luego  $F_m$  puede calcularse, y compararse con el valor teórico de la distribución  $F$ .

Ejemplo En un laboratorio dos máquinas A y B llenan con el producto XYZ, sendos envases. Se afirma que ambas son igualmente precisas en el volumen que llenan. Para probarlo se toman muestras de tamaño 28 de A y 21 de B obteniendo varianzas muestrales de valor 100 y 40 respectivamente. Si las varianzas teóricas son iguales  $F_m = 2.5$  y  $PROB F(2.5, 27, 20) = 0.98$ . Luego, el valor  $F_m = 2.5$  está entonces ubicado en un extremo de la distribución ( $p = 0.98$ ) y es poco probable, por lo que se rechaza la hipótesis.

## 2.2 Estimación por intervalos

Un estimador puntual proporciona un valor aproximado de un parámetro  $\mu$  de una población, pero hay situaciones donde esa información es insuficiente y se desea responder a la pregunta: ¿cuán próxima es la estimación al parámetro? La estimación por intervalos de confianza da respuesta a la interrogante anterior.

Sea el conjunto  $\{X_1; X_2; \dots; X_n\}$  una Muestra Aleatoria Simple (MA S), de una Variable Aleatoria  $X$ , y el parámetro  $\mu$  el que se desea estimar. Para dar una respuesta en términos probabilísticos se construye un intervalo de confianza  $[\inf; \sup]$  que con un nivel de confianza previamente especificado contenga al parámetro  $\mu$ . Más específicamente:

Definición 18 Para establecer un intervalo de confianza  $[\inf; \sup]$ , de un parámetro  $\mu$ , donde  $1 - \alpha$  es el nivel de confianza deseado, se hallarán dos estadísticos  $\inf$  y  $\sup$  tal que la probabilidad de que el intervalo cubra a  $\mu$  es  $1 - \alpha$ .

$$Prob(\inf \leq X_1; \dots; X_n \leq \sup) = 1 - \alpha \quad (2.12)$$

¿Cómo interpretar este concepto en términos frecuenciales?

Si se toman 100 muestras aleatorias de la población, y se construye un intervalo de confianza  $[\inf; \sup]$  del 95 % para cada una de ellas, se espera que en 95 de ellas  $\mu$  se encuentre entre  $\inf$  y  $\sup$ , o sea que el intervalo  $[\inf; \sup]$  contenga al parámetro. Obsérvese con detenimiento las siguientes consideraciones:

<sup>2</sup> que lo aleatorio es el intervalo  $[\inf; \sup]$  que se modifica para cada muestra mientras que el parámetro  $\mu$  es un desconocido pero...jo

<sup>2</sup> que es el intervalo aleatorio a priori, creado con los estadísticos inf y sup, el que tiene el nivel de confianza pero una vez observados los valores de la muestra y estimados inf y sup a partir de los valores muestrales, ya este intervalo concreto no puede interpretarse en términos de probabilidades, y contendrá o no al parámetro

## 2.2.1 Intervalos de la media y de la varianza

Estimación de la media <sup>1</sup> cuando se conoce <sup>3/4</sup>

Se determinará el intervalo de confianza de probabilidad  $(1 - \alpha)$  del valor esperado <sup>1</sup> de la variable aleatoria  $X$ ,  $N(\mu; \sigma^2)$ , cuando se conoce <sup>3/4</sup>.

Si  $X$  es  $N(\mu; \sigma^2)$  y la muestra de tamaño  $n$ , su media muestral  $m$  se distribuye también normalmente pero  $N(\mu; \frac{\sigma^2}{n})$ . Por lo tanto

$$PROBIT\left(\frac{\alpha}{2}\right) \cdot \frac{m - \mu}{\frac{\sigma}{\sqrt{n}}} = PROBIT\left(1 - \frac{\alpha}{2}\right) \quad (2.13)$$

Luego el intervalo de confianza está dado por:

$$m \pm PROBIT\left(1 - \frac{\alpha}{2}\right) \frac{\sigma}{\sqrt{n}} = \mu \pm PROBIT\left(1 - \frac{\alpha}{2}\right) \frac{\sigma}{\sqrt{n}} \quad (2.14)$$

Una de las aplicaciones principales de la distribución  $\chi^2$  es la estimación de la varianza <sup>3/4</sup> de la población.

Estimación de la varianza <sup>3/4</sup>

Un intervalo de confianza de probabilidad  $(1 - \alpha)$  de la varianza <sup>3/4</sup> es:

$$\left[ \frac{(n-1)S^2}{\chi^2_{\alpha/2}(n-1)} ; \frac{(n-1)S^2}{\chi^2_{1-\alpha/2}(n-1)} \right]$$

Ejercicio Se toma una muestra del número de empleados de llaves transportados en taxi diariamente. Los resultados son:

189, 231, 178, 220, 220, 213, 195, 170, 205, 223, 189, 198, 210

199, 211, 198, 210, 222, 203, 198, 194, 205, 203, 199, 205, 211

Halle un intervalo de un 88% de confianza de <sup>3/4</sup> la desviación estándar.

## Estimación de la media cuando la varianza es desconocida

Cuando se desea estimar la media  $\mu$  a partir de una muestra de tamaño  $n$  y se desconoce  $\sigma^2$ , se debe recurrir a la distribución  $t$  de Student con  $n - 1$  grados de libertad. Si  $\bar{m}$  y  $s$  son la media y la desviación estándar muestral respectivamente de una muestra  $x_1, \dots, x_n$  de una variable  $X \sim N(\mu, \sigma^2)$  entonces de las ecuaciones 2.8, 2.9 se deduce que

Proposición 19  $t = \frac{\bar{m} - \mu}{\frac{s}{\sqrt{n}}}$  se distribuye como una  $t$  de Student con  $n - 1$  grados de libertad

La distribución acumulada en SAS de la distribución  $t_{g,l}$  es  $PROBIT(x, g)$  y su inversa  $TINV(p, g)$ .

Su aplicación inmediata es que

$$\bar{m} \pm z_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \sim TINV\left(1 - \frac{\alpha}{2}; n - 1\right) \quad (2.15)$$

es un intervalo de confianza de probabilidad  $(1 - \alpha)$  del valor esperado  $\mu$ .

Ejercicio halle un intervalo de confianza del 90% del valor esperado del número de pasajeros diarios del ejercicio anterior.

### 2.2.2 Tamaño muestral

Cuando se estima  $\mu$  y  $\sigma^2$  es conocido, el intervalo de confianza tiene una longitud de

$$L = 2 z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{N}} \sim PROBIT\left(1 - \frac{\alpha}{2}\right)$$

estableciéndose una relación simple entre la confianza, la precisión  $L$  del resultado (longitud del intervalo), y el tamaño de la muestra. También se considera en algunos textos el error  $E = L/2$ ; obsérvese que no se puede fijar arbitrariamente las 3 cantidades. Si deseo una cierta precisión con un grado de confianza  $1 - \alpha$ , entonces el tamaño de la muestra queda determinado por la fórmula anterior. Más precisamente

$$N = \text{CEIL} \left[ \left( \frac{2 z_{\frac{\alpha}{2}} \text{PROBIT}\left(1 - \frac{\alpha}{2}\right)}{L} \right)^2 \right] \quad (2.16)$$

obsérvese que la función  $\text{ceil}$  es para obtener el entero más próximo por exceso

Si  $\sigma^2$  es desconocido, en la ecuación 2.16 debería cambiarse  $\sigma^2$  por  $S^2$ , y  $\text{PROBIT}$  por  $T$  III V, pero ésta depende de los grados de libertad  $N - 1$ , además, como la varianza muestral  $S^2$  se obtiene del muestreo no se puede conocer a priori. Es por esta razón que para este caso se necesita:

- una estimación anterior de  $\sigma^2$  con su intervalo de confianza adecuado
- un  $N$  razonablemente grande como para sustituir la distribución  $T$  por la normal.

Si no se dispone de una estimación anterior de  $\sigma^2$ , es necesario hacer un muestreo previo pequeño para estimarla (muestreo preliminar):

- Calcule un intervalo de confianza de  $\sigma^2$ , (con confianza = 0.8, por ejemplo), para usar su extremo superior como  $\sigma^2$  si desea ser cauteloso aunque el  $N$  obtenido sea grande
- Úsese directamente la desviación estándar muestral estimada si quiere ser pragmático y confía en la estimación preliminar, corriendo el riesgo de que el  $N$  sea subestimado y el intervalo ...nal supere la precisión deseada.

### 2.2.3 Intervalo de confianza para proporciones

La información de que se dispone para estimar una proporción,  $(p)$ , es el número  $x$  de veces que el evento considerado ocurre en  $n$  ensayos (tamaño muestral). Como es una distribución binomial el valor esperado de  $x$  es  $np$ , el estadístico  $\frac{x}{n}$ , (frecuencia relativa), es un estimador insesgado de  $p$  ( $E(\frac{x}{n}) = p$ ). La varianza de este estimador  $\text{Var}(\frac{x}{n}) = \frac{p(1-p)}{n}$   $\neq 0$  si  $n \neq 1$ ; luego  $\frac{x}{n}$  es además consistente.

Usando la distribución binomial, para cada  $n$  es posible construir intervalos exactos de confianza de cada proporción  $p$ , pero esto es algo complicado. Es por ello que en general se intenta construir intervalos aproximados por un método más simple usando la distribución normal.

Si se cumple la regla práctica para aproximar una binomial con una normal:  $np > 5$   $n(1 - p) > 5$  podemos construir intervalos aproximados de una proporción asimilándolos a los de una normal con igual media y varianza.

El intervalo de confianza aproximado es:

$$\frac{x}{n} \pm \sqrt{\frac{p(1-p)}{n}} \text{PROBIT}\left(1 - \frac{\alpha}{2}\right) \quad (2.17)$$

Como  $p$  es desconocido se tienen dos opciones:

<sup>2</sup> Sustituir  $p(1 \pm p)$  por  $\frac{1}{4}$  en el intervalo anterior de modo de introducir el valor máximo que el producto  $p(1 \pm p)$  puede tener, y así obtener un intervalo más grande pero seguro.

<sup>2</sup> Usar como valor de  $p$  su estimador  $\frac{x}{n}$ , criterio este más pragmático para definir el intervalo pero menos confiable.

Cuando  $p$  es próximo a  $1/2$  los intervalos son semejantes. En cuanto al tamaño de la muestra para obtener una precisión  $L$  con confianza  $(1 - \alpha)$ , de la ecuación 2.17:

$$N = 4pe(1 \pm pe) \frac{\mu_{\text{PROBIT}}(1 \pm \frac{\alpha}{2})^2}{L^2} \quad (2.18)$$

donde  $pe$  es una estimación previa de  $p$ . Si se quiere una estimación conservadora pero segura, o en caso de no disponer de  $pe$ , sustituir  $pe$  por  $1/2$  en la ecuación 2.18 para obtener una estimación por exceso pero segura del tamaño muestral.

$$N = \frac{\mu_{\text{PROBIT}}(1 \pm \frac{\alpha}{2})^2}{L^2} \quad (2.19)$$

## 2.3 Prueba de hipótesis

Una hipótesis es una afirmación acerca de una característica desconocida de una población. Por ejemplo  $\mu = 123$ ; si existe suficiente evidencia experimental que apoye la hipótesis.

A pesar de que el significado real de la hipótesis se refiere a características significativas de la población, (al valor esperado, a la desviación estándar, al máximo, etc), toda hipótesis en un test se establece en términos de parámetros de una variable aleatoria. Por fortuna, para la mayoría de las distribuciones estos parámetros coinciden con el valor esperado, la varianza, etc. Así cuando se decide entre dos hipótesis lo que realmente se hace es decidir entre posibles valores de los parámetros. Más precisamente, se opta entre valores de los parámetros para un mismo tipo de distribución. Por ejemplo, si una población es normal  $N(20; 3/4)$  o normal  $N(23; 3/4)$ :

Ilustrémoslo con el siguiente ejemplo: disponemos de una empaquetadora  $XX$  que llena 120 bdsas por minuto, con una desviación estándar  $3/4 = 5$ . Se nos ofrece comprar otra  $LL$  de la cual se afirma que su media es superior, y que someteremos a un test: o se invertirá dinero en comprar  $LL$  salvo que exista clara evidencia experimental que indique que  $LL$  tiene una media mayor. Se supondrá que ambas máquinas tienen la misma desviación

estándar conocida. Se nos permite hacer una prueba sobre  $\mu$  de modo de tomar una muestra de varios minutos de producción  $X_1, \dots, X_n$ .

La hipótesis nula, la hipótesis que no deseamos abandonar salvo que exista suficiente evidencia en su contra, es que la media de  $\mu$  es igual a la de  $X$ . La hipótesis alternativa es que la media es mayor.

### 2.3.1 Conceptos básicos

Se formulará este problema en general:

$H_0: \mu = \mu_0 = 120$  Prueba unilateral

$H_a: \mu = \mu_a > \mu_0$   $\frac{3}{4}$  conocida

Cuando la hipótesis alternativa es que la media es diferente pero no importa si es mayor o menor se formula un test bilateral:

$H_0: \mu = \mu_0 = 120$  Prueba bilateral

$H_a: \mu = \mu_a \neq \mu_0$   $\frac{3}{4}$  conocida

Se toma la muestra aleatoria  $X_1, \dots, X_n$ , del número de bdsas producido por minuto. Se calcula la media muestral  $\bar{m}$  y se compara con un número  $L$  preestablecido, que más adelante se determinará:

#### CRITERIO DE DECISIÓN

R egión de aceptación de $H_0$	$\bar{m} < L$	R egión de rechazo de $H_0$	$\bar{m} > L$
OPTAR POR $H_0$		OPTAR POR $H_a$	

Entre la realidad y la decisión tomada, pueden establecerse 4 situaciones diferentes, 2 aciertos y 2 errores que se reflejan en el siguiente cuadro

	Decido $H_0$	Decido $H_a$
Realidad $H_0$		error de primer tipo <sup>®</sup>
Realidad $H_a$	error de segundo tipo <sup>-</sup>	Potencia = $1 - \beta$

Lo deseable sería reducir al mínimo la probabilidad de los errores pero para una muestra de tamaño  $n$  esto no es posible. De los dos errores indicados, el más importante, comprar la nueva máquina siendo su media igual a la que tengo, se le designa como error de primer tipo o nivel de significación de la prueba<sup>®</sup>: Es la probabilidad de optar por  $H_a$  siendo verdadera la  $H_0$ .<sup>®</sup> =  $\text{pr}ob(H_a = H_0)$

Es el error que no deseo cometer y lo controlo asignándole una probabilidad baja<sup>®</sup> (nivel de significación). Fijado el nivel de significación, es

deseable que la potencia del test sea alta o lo que es equivalente que el error de segundo tipo sea bajo ( $Potencia = 1 - \beta$ ).

Error de tipo II  $\beta = \text{prdb}(H_0 = H_a)$ :

Potencia del test  $= \text{prdb}(H_a = H_a)$ .

Para establecer el criterio de decisión se determina  $L$  de manera que  $\text{Prdb}(m > L = H_0) = \alpha$ .

Fijado  $L$ , el error de segundo es tipo  $\text{Prdb}(m < L = H_a) = \beta$

En términos jurídicos, el error de primer tipo es condenar al inocente, y el de segundo, absolver al culpable.

Retomando el ejemplo, como la desviación estándar es conocida, por las propiedades de la media muestral,  $L$  está dado por:

$$L = \mu_0 + \frac{\sigma}{\sqrt{n}} \text{PROBIT}(\alpha) \quad (2.20)$$

y el error de segundo tipo

$$\beta = \text{PROBNORM}\left(\frac{L - \mu_a}{\frac{\sigma}{\sqrt{n}}}\right) \quad (2.21)$$

Tamaño muestral en función de  $\alpha$  y  $\beta$

En la página anterior se fijaron el error de primer tipo y el tamaño de la muestra, en consecuencia, el error de segundo tipo quedó determinado por  $\alpha$  y  $n$ . Si, por el contrario, se desea controlar a priori ambos errores y se fijan  $\alpha$  y  $\beta$ , el tamaño adecuado de la muestra  $n$  se puede calcular a partir de  $\alpha$  y  $\beta$ .

Se incluye a continuación los resultados para los 3 posibles test, se observará que la alternativa  $\mu_a$ , debe fijarse en un valor preestablecido

Caso Hipótesis simple contra alternativa simple con desviación estándar conocida

I)

$H_0: \mu \geq \mu_0$  Prueba unilateral,  $\sigma$  conocida

$H_a: \mu < \mu_0$

$$N = \frac{\sigma^2 [\text{PROBIT}(\alpha) + \text{PROBIT}(\beta)]^2}{(\mu_0 - \mu_a)^2}$$

II) Si la alternativa es una prueba bilateral,  $\mu_a \notin \mu_0$ , el paréntesis del numerador debe sustituirse por:  $\text{PROBIT}(\alpha/2) + \text{PROBIT}(\beta)$ , pero el valor de  $N$  es aproximado



### 2.3.2 Como programar las pruebas

En la formulación anterior se tomó la decisión comparando la media muestral  $\bar{m}$  con el valor  $\mu_0$ . Sin embargo, lo habitual es estandarizar el resultado como se indica en los test siguientes:

Prueba sobre la media (con desviación estándar  $\sigma$  conocida).

Estadístico de prueba

$$Est = \left( \frac{\bar{m} - \mu_0}{\frac{\sigma}{\sqrt{N}}} \right)$$

H hipótesis nula  $\mu = \mu_0$

Prueba sobre la media  $\mu$  conocida

H ip.alternativa

A ccepta H o

valor p

$\mu_a > \mu_0$

$Est < PROBIT(1-\alpha)$

$1-PROBNORM(Est)$

$\mu_a < \mu_0$

$Est > PROBIT(\alpha)$

$PROBNORM(Est)$

$\mu_a \neq \mu_0$

$PROBIT(\frac{\alpha}{2}) < Est < PROBIT(1-\frac{\alpha}{2})$

$2(1-PROBNORM(|Est|))$

Prueba sobre la media (con desviación estándar desconocida).

Estadístico de prueba

$$Est = \left( \frac{\bar{m} - \mu_0}{\frac{s}{\sqrt{N}}} \right)$$

H hipótesis nula  $\mu = \mu_0$

Prueba sobre la media  $\mu$  desconocida

H ip.alternativa

A ccepta H o

valor p

$\mu_a > \mu_0$

$Est < TINV(1-\alpha; N-1)$

$1-PROBT(Est; N-1)$

$\mu_a < \mu_0$

$Est > TINV(\alpha; N-1)$

$PROBT(Est; N-1)$

$\mu_a \neq \mu_0$

$PROBT(\frac{\alpha}{2}; N-1) < Est < PROBT(1-\frac{\alpha}{2}; N-1)$

$2(1-PROBT(|Est|; N-1))$

Prueba sobre la varianza

Estadístico de prueba

$$Est = (j-1) \frac{S^2}{\sigma_0^2}$$

H hipótesis nula  $\mu_a = \mu_o$

prueba acerca de la varianza

H ip.alternativa

A ccepta H o

valor p

$\mu_a > \mu_o$

Est < CIII V (1- $\alpha$ , N -1)

1-PROBCHI(Est/N -1)

$\mu_a < \mu_o$

Est > CIII V ( $\alpha$ , N -1)

PROBCHI(Est/N -1)

$\mu_a \neq \mu_o$

CIII V ( $\frac{\alpha}{2}$ , N -1) < Est < CIII V (1- $\frac{\alpha}{2}$ , N -1)

2(1-PROBCHI(jEstj, N -1))

Valor p

Es posible que el valor del estadístico (Est) sea próximo al límite de decisión y se quiere expresar esta proximidad en términos probabilísticos. Si se desea saber cuán cerca se está de la probabilidad  $\alpha$ , se determina el valor de probabilidad, valor p, que es:

El mínimo nivel de signi...cación para el cual los datos observados indican que se tendría que rechazar la hipótesis nula.

Esta probabilidad se calcula de manera idéntica a la forma en que se asignó la región de rechazo de modo que si la media muestral coincide con el o los límites de la región el valor p es  $\alpha$ .

En el caso anterior,  $p = 1 - \text{PROBCHI}(\text{Est})$  y se le compara con el error de primer tipo. Por ejemplo, si se diseñó un test con  $\alpha = 0.05$  y el valor p es 0.06 U d. acepta la hipótesis nula pero le queda la duda ya que con otra muestra similar podría haber dado un valor  $p = 0.047$  y optar por la alternativa. Sin embargo, si  $p = 0.30$  U d. aceptaría la hipótesis nula con mayor convicción.

### 2.3.3 Pruebas para dos poblaciones

Prueba sobre la diferencia de medias

Sean  $x_1, \dots, x_n$  una muestra de una variable  $X \sim N(\mu_x, \sigma_x^2)$  y  $y_1, \dots, y_n$  una muestra de una variable  $Y \sim N(\mu_y, \sigma_y^2)$ .

CA S0 : Varianzas de ambas poblaciones conocidas.

Estadístico de prueba

$$\text{Est} = \left( \frac{\mu_y - \mu_x}{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}} \right) \quad (2.22)$$

H hipótesis nula  $\mu_x = \mu_y$

### Prueba de igualdad de medias con varianzas conocidas

Hipótesis alternativa	Acepta H <sub>0</sub>	valor p
$\mu_y > \mu_x$	$\text{Est} < \text{PROBIT}(1 - \alpha)$	$1 - \text{PROBNORM}(\text{Est})$
$\mu_y \neq \mu_x$	$\text{PROBIT}(\frac{\alpha}{2}) < \text{Est} < \text{PROBIT}(1 - \frac{\alpha}{2})$	$2(1 - \text{PROBNORM}( \text{Est} ))$

CA S0: Varianzas de ambas poblaciones iguales pero desconocidas.

Supóngase que las varianzas son iguales pero que no se conoce su valor:

$$\sigma_{\mu_x}^2 = \sigma_{\mu_y}^2 = \sigma^2 \text{ (desconocida).}$$

Estadístico de prueba

$$\text{Est} = \left( \frac{\bar{m}_y - \bar{m}_x}{S \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} \right) \quad (2.23)$$

Donde  $S^2$  es un promedio de las varianzas muestrales ponderado por el tamaño muestral:

$$S^2 = \frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}$$

Sea  $n_z = n_x + n_y - 2$  la suma de los grados de libertad asociados a cada muestra

Hipótesis nula  $\mu_x = \mu_y$

Prueba de igualdad de medias con varianzas desconocidas iguales

Hipótesis alternativa	Acepta H <sub>0</sub>	valor p
$\mu_y > \mu_x$	$\text{Est} < \text{TINV}(\alpha; n_z)$	$1 - \text{PROBT}(\text{Est}; n_z)$
$\mu_y \neq \mu_x$	$\text{TINV}(\frac{\alpha}{2}; n_z) < \text{Est} < \text{TINV}(1 - \frac{\alpha}{2}; n_z)$	$2(1 - \text{PROBT}( \text{Est} ; n_z))$

### Pruebas sobre observaciones apareadas

Cuando las observaciones de las dos variables se hacen sobre las mismas unidades experimentales

y en consecuencia los tamaños muestrales son iguales, debe considerarse una nueva variable

$Z = Y - X$ , y reducir el análisis de la diferencias de medias al estudio de si la media de la variable  $Z$  es 0. Es decir, en el caso de observaciones apareadas, las pruebas sobre dos poblaciones se reducen a las pruebas de una población si se considera como nueva variable la diferencia de las variables que se desea comparar.

Prueba sobre la igualdad de varianzas

Sean  $\{X_1, \dots, X_{n_x}\}$  una muestra de una variable  $X \sim (1, \sigma_x^2)$  y  $\{Y_1, \dots, Y_{n_y}\}$  una muestra de una variable  $Y \sim (1, \sigma_y^2)$ .

Estadístico de prueba

$$F_{\text{est}} = \frac{S_y^2}{S_x^2} \quad (2.24)$$

Hipótesis nula  $\sigma_x = \sigma_y = \sigma$  (desconocida)

Prueba de igualdad de varianzas

Hipótesis alternativa

Acepta  $H_0$

valor p

$$\sigma_y > \sigma_x$$

$$F_{\text{est}} < F_{1-\alpha/2, n_y-1, n_x-1}$$

$$1 - \text{PROB F}(F_{\text{est}}, n_y-1, n_x-1)$$

$$\sigma_y \neq \sigma_x$$

$$F_{\alpha/2, n_y-1, n_x-1} < F_{\text{est}} < F_{1-\alpha/2, n_y-1, n_x-1}$$

$$2(1 - \text{PROB F}(F_{\text{est}}, n_y-1, n_x-1))$$

Pruebas clásicas para poblaciones normales con el SAS

Prueba sobre medias

Una población

$H_0: \mu = \mu_0$  Prueba unilateral,  $\sigma$  desconocida

$H_a: \mu > \mu_0$  o  $\mu < \mu_0$

Probar que el valor esperado de  $X$  es  $\mu_0$  es lo mismo que probar que el de la variable  $Z = (X - \mu_0)/\sigma$  es 0. Para hacer ese test se toma una muestra  $z_1, z_2, \dots, z_n$  de ella y se aplica el PROC UNIVARIATE que determina la probabilidad de que la media muestral se desvíe de 0 un valor igual al observado. El SAS incluye el valor p.

$H_0: \mu = \mu_0$  Prueba unilateral,  $\sigma$  desconocida

$H_a: \mu > \mu_0$  o  $\mu < \mu_0$

Idéntico al anterior pero observando que el valor p debe ser la mitad de lo que indica el UNIVARIATE y además, la media muestral debe tener el signo adecuado al test.

Dos poblaciones

El procedimiento PROCTTEST permite comparar medias de dos variables  $X$  e  $Y$ , a partir de sendas muestras de tamaño arbitrario, tanto en el caso de que sus varianzas sean iguales o distintas.

Prueba

H<sub>0</sub> →  $\mu_x = \mu_y$

H<sub>a</sub> →  $\mu_x \neq \mu_y$

Para probar si dos variables X e Y tienen el mismo valor esperado se toman muestras  $fX_1; \dots; X_{nx}$ g de X e  $fY_1; \dots; Y_{ny}$ g de Y. El PROC TTEST del SAS permite realizar la prueba de igualdad de medias.

Ejemplo incluya 'ttest.sas' del directorio de ejercicios.

Prueba

H<sub>0</sub> →  $\mu_x = \mu_y$

H<sub>a</sub> →  $\mu_x < \mu_y$  o  $\mu_x > \mu_y$

Idéntico al anterior pero observando que el valor p debe ser la mitad de lo que indica el TTEST y además la media muestral de Y mayor que la de X si  $\mu_x < \mu_y$ .

Prueba

H<sub>0</sub> →  $\mu_{x+d} = \mu_y$

H<sub>a</sub> →  $\mu_{x+d} < \mu_y$

Para probar si las variables X e Y cumplen que el valor esperado de Y es d unidades mayor que el de X, se toman muestras  $fX_1; \dots; X_{nx}$ g de X e  $fY_1; \dots; Y_{ny}$ g de Y. Se crea una nueva variable  $Z = X + d$  que tendrá por valor esperado el de  $X + d$  y se reduce el problema al Test anterior usando el PROC TTEST del SAS con Z e Y.

### Prueba sobre varianzas

Para probar si dos variables X e Y tienen varianzas iguales se toman muestras  $fX_1; \dots; X_{nx}$ g de X e  $fY_1; \dots; Y_{ny}$ g de Y. El PROC TTEST del SAS permite realizar el test de igualdad de varianzas.

Prueba bilateral:

H<sub>0</sub> →  $\sigma_x^2 = \sigma_y^2$

H<sub>a</sub> →  $\sigma_x^2 \neq \sigma_y^2$

Incluye el valor del estadístico F y el valor p asociado

ooooo

Prueba unilateral:

H<sub>0</sub> →  $\sigma_x^2 = \sigma_y^2$

H<sub>a</sub> →  $\sigma_x^2 > \sigma_y^2$

Si el test es unilateral, debe tomarse como valor p la mitad, y observar si el test tiene el signo coherente con la desigualdad de la hipótesis alternativa, o sea, la varianza muestral de Y debe ser mayor que la de X.

### 2.3.4 Pruebas de ajuste de distribuciones

Prueba de Kolmogorov para el ajuste de distribuciones muestrales a las distribuciones continuas, donde los parámetros estén totalmente identificados.

La siguiente función de Kolmogorov juega un rol fundamental en los tests de ajuste de distribuciones empíricas.

$$K(z) = \sum_{k=-\infty}^{\infty} (j-1)^k \exp(j 2 k^2 z^2)$$

Algunos valores significativos de esta función se incluyen a continuación:

K(z)	.75	.80	.90	.95	.975	.99
z	1.02	1.07	1.22	1.36	1.48	1.64

Dadas dos distribuciones empíricas  $F_n(x)$  y  $G_m(x)$  y una teórica  $F(x)$ , sean:

$$D_n = \max_j |F_n(x_j) - F(x_j)|$$

$$D_{nm} = \max_j |F_n(x_j) - G_m(x_j)|$$

donde desde un punto de vista práctico los máximos se obtienen considerando solamente los valores discretos de las distribuciones muestrales empíricas.

Prueba de ajuste de la distribución empírica  $F_n(x)$  a la teórica  $F(x)$  totalmente especificada.

La hipótesis nula: Los valores muestrales  $X_1, \dots, X_n$  provienen de la distribución teórica  $F(x)$ .

Suposición estadística: para  $n$  "razonablemente" grande se cumple

$$\text{Prob}\left[\sqrt{n} D_n < z\right] = K(z)$$

Diseño de la Prueba

Establecido un nivel de significación  $\alpha$ , este valor debe ser igual a  $1 - K(z)$ . Se halla, entonces, el  $z$  asociado a  $K(z)$ .

Cómo decidir?

Dada una muestra de tamaño  $n$ , se debe hallar  $D_n$  y verificar la ecuación anterior. Si

$$p_{(N)} = D_N \cdot z$$

se acepta la hipótesis nula

Valor p: si se considera el valor de z que satisface la igualdad anterior, entonces el valor p del test está dado por  $1 - K(z)$ :

Prueba de ajuste sobre si ambas distribuciones empíricas  $F_n(x)$  y  $F_m(x)$  provienen de la misma distribución teórica  $F(x)$  desconocida.

Para n y m razonablemente grandes se cumple:

$$Pr\left[\frac{r}{n+m} D_{nm} < z\right] = K(z)$$

Luego dadas sendas muestras de tamaños m y n, se debe hallar  $D_{nm}$  y

aplicar la ecuación anterior. El valor p del test está dado por  $1 - K(z)$ .

### 2.3.5 Prueba de independencia

X e Y son independientes cuando la estructura de probabilidad de X no cambia por los valores que toma Y.

Dadas dos Variables Aleatorias X e Y con valores esperados  $\mu_x$  y  $\mu_y$  y desviaciones standard  $\sigma_x$  y  $\sigma_y$  respectivamente, se define covarianza de X e Y:

$$COV(X; Y) = E[(X - \mu_x)(Y - \mu_y)]$$

y correlación entre X e Y:

$$\rho = CORR(X; Y) = \frac{COV(X; Y)}{\sigma_x \sigma_y}$$

La correlación expresa el grado de asociación lineal entre dos variables y

siempre es un número entre -1 y 1. Si las variables son independientes su correlación es 0.

El PROC CORR del SAS calcula las correlaciones muestrales entre dos o más variables y al mismo tiempo realiza una prueba de hipótesis acerca de si las variables son independientes.