

MUEST. T6. , MUESTREO ESTRATIFICADO :

2. ESTIM. LINEALES, VARIANZAS y SUS ESTIMACIONES.
 4. x. PRINCIPIOS BÁSICOS de la ESTRATIFICACIÓN.
 5. x. CONSTRUCCIÓN y NÚMERO de ESTRATOS.
 3. AFIJACIÓN de la MUESTRA con 1 CARACTERÍSTICA.
 6. GANANCIA de PRECISIÓN
-

1. CONCEPTO de MUESTREO ESTRATIFICADO.

En el muestreo irrestuido aleatorio la población se considera homogénea respecto a la característica observada. Las estimaciones obtenidas son buenas a un coste bajo.

En el muestreo estratificado, la población se considera heterogénea respecto a la característica estudiada, se subdivide en subpoblaciones lo más homogéneas posibles, con lo que se obtienen estimaciones más precisas con el mismo coste.

La población heterogénea con N unidades, $\{U_i\}, i=1 \dots N$ se subdivide (partición) en L estratos de \neq tamaños, $N_1, \dots, N_h, \dots, N_L$

$\{U_{hi}\} \mid \begin{array}{l} h=1 \dots L \rightarrow n^\circ \text{ estratos} \text{ pblar.} \\ i=1 \dots N_h \rightarrow n^\circ \text{ unidades por estrato} \end{array}$

De cada estrato se selecciona una muestra aleatoria, el tipo de muestreo puede ser \neq para cada estrato, de manera independiente.

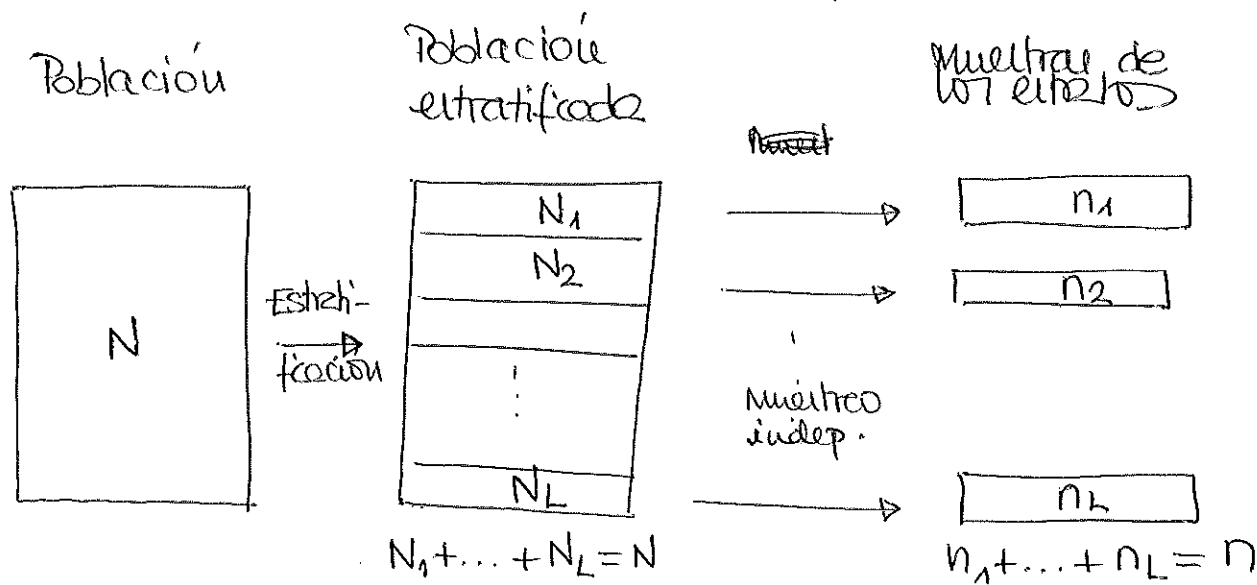
$\{U_{hi}\} \mid \begin{array}{l} h=1 \dots L \rightarrow n^\circ \text{ estratos} \\ i=1 \dots n_h \rightarrow n^\circ \text{ unidades muestrales por estrato} \end{array}$

La homogeneidad / heterogeneidad se estudia a través de la varianza \cong precisión del valor promedio.

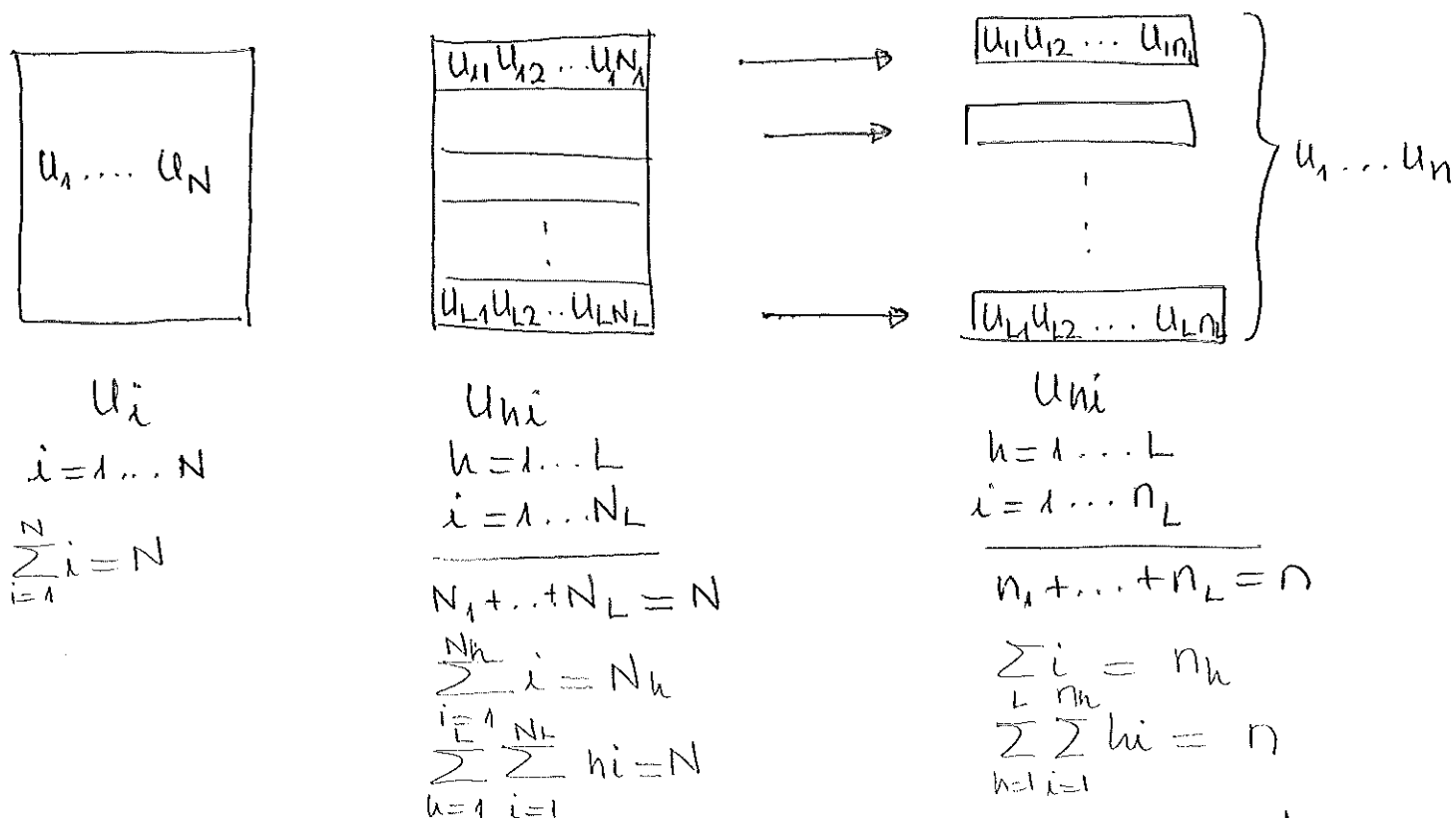
Homogénea \cong varianzas iguales

Heterogénea \cong varianzas distintas

El procedimiento del muestreo estratificado se puede resumir mediante el siguiente esquema:



que se puede detallar más especificando las unidades en cada grupo



Nótese que la estratificación es a nivel poblacional, y que en la muestra están presentes todos los estratos

El muestreo estratificado es aconsejable por los siguientes motivos:

- Permite obtener estimaciones para cada subpoblación.
- Puede generar ganancia en precisión.
Al dividir la población heterogénea en estratos homogéneos, puede que la variación en cada estrato sea menor que la variación en toda la población.
- Se pueden utilizar distintos tipos de muestreo para cada estrato, lo que permite reducir el coste.
- Si existe una variable precisa para la estratificación, correlacionada con la variable, se permite dividir la población en estratos homogéneos.

2 - ESTIMADORES LINEALES, VARIANZAS Y SUS ESTIMACIONES

Aunque se pueden utilizar distintos tipos de muestreo para cada estrato, vamos a suponer por comodidad que el muestreo se realiza de manera independiente y utilizando muestreo aleatorio con probabilidades iguales en cada estrato. Distinguiamos los casos de SIN y CON reposición.

Población:

Dividimos la población en L estratos. Sobre cada unidad poblacional observamos la característica X . El parámetro a estimar es el valor poblacional de la característica.

$$\theta = \sum_{i=1}^N Y_i \longrightarrow \theta = \sum_{h=1}^L \sum_{i=1}^{(n_h)} Y_{hi} \quad \left. \begin{array}{l} h \rightarrow \text{estrato} \\ i \rightarrow \text{unidad dentro estrato} \end{array} \right\}$$

ESTIMADORES LINEALES:

Utilizamos como estimador lineal

$$\hat{\theta} = \sum_{h=1}^L \sum_{i=1}^{(n_h)} w_{hi} Y_{hi}$$

donde w_{hi} es el peso de muestreo de Y_{hi} , y representa el u^0 de unidades elementales de la unidad compuesta o su ponderación o importancia.

Para comprobar su insesgadez, usaremos de la variable auxiliar e_{hi} , $e_{hi} \equiv u^0$ de veces que aparece u_{hi} en la muestra del estrato h .

Si reposición:

Utilizaremos como estimador insesgado la suma extendida a todos los extratos de los estimadores unidades insesgadas de Horvitz y Thompson en cada extrato.

$$(SR) \rightarrow \hat{\theta}_{st} = \sum_{h=1}^L \sum_{i=1}^{N_h} \frac{y_{hi}}{\pi_{hi}} \quad , \text{ donde } \pi_{hi} \equiv \text{probab. de que } u_{hi} \in S_h .$$

Para demostrar la insesgadez, acortemos a e_{hi} con probab. $\pi_{hi} = \frac{n_h}{N_h}$ (probab. iguales).

$$e_{hi} = \begin{cases} 1 & \text{si } u_{hi} \in S_h \\ 0 & \text{si } u_{hi} \notin S_h \end{cases} \quad \text{con probab. } 1 - \pi_{hi} = 1 - \frac{n_h}{N_h}$$

$$e_{hi} \rightarrow \text{Ber}(\pi_{hi}) \text{ con } E[e_{hi}] = \pi_{hi} = \frac{n_h}{N_h}$$

$$\begin{aligned} E[\hat{\theta}] &= E\left[\sum_{h=1}^L \sum_{i=1}^{N_h} y_{hi} w_{hi}\right] = E\left[\sum_{h=1}^L \sum_{i=1}^{N_h} y_{hi} w_{hi} e_{hi}\right] = \\ &= \sum_{h=1}^L \sum_{i=1}^{N_h} y_{hi} w_{hi} \underbrace{E[e_{hi}]}_{\pi_{hi}} = \theta \Leftrightarrow w_{hi} = \frac{1}{\pi_{hi}}, \forall h, i. \end{aligned}$$

La expresión del estimador se puede particularizar para los parámetros poblacionales más comunes:

$$\text{Total poblacional: } \theta = X = \sum_{h,i} X_{hi} \rightarrow \hat{X}_{st} = \sum_{h=1}^L \hat{X}_h$$

El estimador del total es la suma de los estimadores del total en cada extrato.

$$\begin{aligned} \hat{X}_{st} &= \sum_{h=1}^L \sum_{i=1}^{n_h} \frac{X_{hi}}{\pi_{hi}} \stackrel{\text{p.i.}}{=} \sum_{h=1}^L \sum_{i=1}^{n_h} \frac{X_{hi}}{n_h/N_h} = \sum_{h=1}^L N_h \cdot \underbrace{\frac{1}{n_h} \sum_{i=1}^{n_h} X_{hi}}_{\hat{X}_h = \bar{X}_h} = \\ &= \sum_{h=1}^L N_h \cdot \bar{X}_h = \sum_{h=1}^L \hat{X}_h \quad \Leftrightarrow \hat{X}_h = N_h \bar{X}_h \end{aligned}$$

Media poblacional: Media ponderada de los estimadores de la media en cada extrato, siendo los coef. de ponderación los que marcan la importancia relativa de cada extrato en la poble.

$$\Theta = \bar{X} \Rightarrow y_{hi} = \frac{X_{hi}}{N} \longrightarrow \hat{\bar{X}}_{st} = \sum_{h=1}^L w_h \bar{x}_h \quad / w_h = \frac{N_h}{N}$$

$$\begin{aligned} \Theta = \bar{X} \Rightarrow \hat{\bar{X}}_{st} &= \sum_{h=1}^L \frac{1}{N} \sum_{i=1}^{n_h} \frac{X_{hi}}{\pi_{hi}} = \sum_{h=1}^L \frac{1}{N} \sum_{i=1}^{n_h} \frac{X_{hi}}{n_h/N_h} = \\ &= \sum_{h=1}^L \frac{N_h}{N} \underbrace{\sum_{i=1}^{n_h} \frac{X_{hi}}{n_h}}_{\bar{x}_h} = \sum_{h=1}^L w_h \cdot \bar{x}_h \quad / \sum w_h = 1. \end{aligned}$$

Total de clase: Suma de los estimadores del total de clase en cada estrato.

$$\begin{aligned} \Theta = A \Rightarrow y_{hi} &= \frac{A_{hi}}{N} \Rightarrow \hat{A}_{st} = \sum_{h=1}^L \hat{A}_h \\ \Theta = A \Rightarrow \hat{A}_{st} &= \sum_{h=1}^L \frac{1}{N} \sum_{i=1}^{n_h} \frac{A_{hi}}{\pi_{hi}} = \sum_{h=1}^L \frac{1}{N} \sum_{i=1}^{n_h} \frac{A_{hi}}{n_h/N_h} = \\ &= \sum_{h=1}^L \frac{N_h}{N} \underbrace{\sum_{i=1}^{n_h} \frac{A_{hi}}{n_h}}_{\hat{P}_h} = \sum_{h=1}^L N_h \cdot \hat{P}_h = \sum_{h=1}^L \hat{A}_h \end{aligned}$$

Proporción poblacional: Media ponderada de los estimadores de la proporción en cada estrato, siendo los coef. de ponderación $w_h = \frac{N_h}{N}$ de serie unitaria.

$$\begin{aligned} \Theta = P \Rightarrow y_{hi} &= \frac{A_{hi}}{N} \Rightarrow \hat{P}_{st} = \sum_{h=1}^L w_h \hat{P}_h \\ \Theta = P \Rightarrow \hat{P}_{st} &= \sum_{h=1}^L \frac{1}{N} \sum_{i=1}^{n_h} \frac{A_{hi}}{\pi_{hi}} = \sum_{h=1}^L \frac{1}{N} \sum_{i=1}^{n_h} \frac{A_{hi}}{n_h/N_h} = \sum_{h=1}^L \frac{N_h}{N} \underbrace{\sum_{i=1}^{n_h} \frac{A_{hi}}{n_h}}_{\hat{P}_h} \\ &= \sum_{h=1}^L \frac{N_h}{N} \hat{P}_h = \sum_{h=1}^L w_h \hat{P}_h \end{aligned}$$

Con reposición:

La muestra estratificada de tamaño n se obtiene seleccionando n_h elementos ($h=1, \dots, L$) de cada uno de los estratos de forma independiente, utilizando muestreo aleatorio simple con reposición.

Para el parámetro poblacional $\Theta = \sum_{h=1}^L \sum_{i=1}^{N_h} Y_{hi}$, se utiliza como estimador la suma extendida a todos los estratos de los estimadores lineales insesgados de Hansen y Hurwitz en cada estrato.

$$\hat{\Theta} = \sum_{h=1}^L \sum_{i=1}^{n_h} \frac{Y_{hi}}{n_h P_{hi}}$$

p. iguales
/ $P_{hi} = \frac{1}{N_h}$

$P_{hi} \equiv$ probab. unitaria de que U_{hi} esté en la muestra del estrato h .

Al ser $P_{hi} = \frac{1}{N_h}$, se observa que $n_h P_{hi} = n_h / N_h = \pi_{hi}$,

por lo que las expresiones del estimador con reposición y sin reposición coinciden, $\hat{\Theta}_{HT} = \hat{\Theta}_{HH}$, lo que se hace extensivo a los estimadores de los parámetros más comunes.

Para demostrarlo insesgado, añadimos a la v.a. auxiliar $e_{hi} \equiv$ veces que U_{hi} está en la muestra del estrato h .

$$e_{hi} \rightarrow B(n_h, \frac{1}{N_h}) \quad E[e_{hi}] = \frac{n_h}{N_h}$$

$$\begin{aligned} E[\hat{\Theta}] &= E\left[\sum_{h=1}^L \sum_{i=1}^{n_h} Y_{hi} W_{hi}\right] = E\left[\sum_{h=1}^L \sum_{i=1}^{N_h} Y_{hi} W_{hi} e_{hi}\right] = \\ &= \sum_{h=1}^L \sum_{i=1}^{N_h} Y_{hi} W_{hi} \underbrace{E[e_{hi}]}_{n_h/N_h} = \Theta \Leftrightarrow W_{hi} = \frac{1}{n_h} = \frac{1}{n_h P_{hi}} \end{aligned}$$

VARIANZAS de los ESTIMADORES

Al estudiar una muestra en lugar de toda la población, se produce un error de muestreo, medido con la desviación típica del estimador.

Un estimador será más preciso cuanto menor sea su variación. En el caso de estimadores insesgados, la acuracidad y la precisión son lo mismo, pues el ECM del estimador coincide con su variación.

Como el muestreo de cada estrato se realiza de forma independiente, la variación del estimador es igual a la suma de las variaciones de cada estrato.

$$V(\hat{\theta}_{st}) = \sum_{h=1}^L V(\hat{\theta}_h)$$

Aunque las expresiones de los estimadores de Horvitz y Thompson coincidirían con los estimadores de Hansen y Hurvitz, no ocurre así con sus variaciones, por lo que se hace necesario distinguir los casos.

SIN reposición: Directamente para los parámetros

$$V(\hat{X}_{st}) = V\left(\sum_{h=1}^L \hat{X}_h\right) \stackrel{\text{indep.}}{=} \sum_{h=1}^L V(\hat{X}_h) = \sum_{h=1}^L N_h^2 (1-f_h) \cdot \frac{S_h^2}{n_h} \quad (\text{mirar +2})$$

$$V(\hat{\bar{X}}_{st}) = V\left(\sum_{h=1}^L W_h \hat{\bar{X}}_h\right) = \sum_{h=1}^L W_h^2 V(\hat{\bar{X}}_h) = \sum_{h=1}^L W_h^2 \cdot (1-f_h) \frac{S_h^2}{n_h} \quad W_h = \frac{N_h}{N}$$

$$V(\hat{A}_{st}) = V\left(\sum_{h=1}^L \hat{A}_h\right) = \sum_{h=1}^L V(\hat{A}_h) = \sum_{h=1}^L N_h^2 \cdot (1-f_h) \cdot \left(\frac{N_h}{N_h-1}\right) \cdot \frac{P_h Q_h}{n_h}$$

$$V(\hat{P}_{st}) = V\left(\sum_{h=1}^L W_h \hat{P}_h\right) = \sum_{h=1}^L W_h^2 V(\hat{P}_h) = \sum_{h=1}^L W_h^2 \cdot (1-f_h) \cdot \left(\frac{N_h}{N_h-1}\right) \frac{P_h Q_h}{n_h}$$

donde:

$$f_h = \frac{n_h}{N_h}$$

$S_h^2 \equiv$ varianza poblacional estrato h

$P_h Q_h \equiv$ varianzas poblac. poblac. dicotómica estrato h .

$$(N_h-1)S_h^2 = N_h \sigma_h^2$$

CON reposición: Tb. directamente

$$V(\hat{X}_{st}) = V\left(\sum_{h=1}^L \hat{X}_h\right) = \sum_{h=1}^L V(\hat{X}_h) = \sum_{h=1}^L N_h^2 \cdot \frac{\sigma_h^2}{n_h}$$

$$V(\hat{\bar{X}}_{st}) = V\left(\sum_{h=1}^L W_h \hat{X}_h\right) = \sum_{h=1}^L W_h^2 V(\hat{X}_h) = \sum_{h=1}^L W_h^2 \cdot \frac{\sigma_h^2}{n_h}$$

$$V(\hat{A}_{st}) = V\left(\sum_{h=1}^L \hat{A}_h\right) = \sum_{h=1}^L V(\hat{A}_h) = \sum_{h=1}^L N_h^2 \cdot \frac{P_h Q_h}{n_h}$$

$$V(\hat{P}_{st}) = V\left(\sum_{h=1}^L W_h \hat{P}_h\right) = \sum_{h=1}^L W_h^2 V(\hat{P}_h) = \sum_{h=1}^L W_h^2 \cdot \frac{P_h Q_h}{n_h}$$

ESTIMACIÓN de las VARIANZAS

Los parámetros poblacionales de los que dependen las varianzas de los estimadores (S_h^2 , $P_h Q_h$) suelen ser desconocidos, por lo que es preciso estimarlos.

Para características cuantitativas, la cuasivarianza muestral de cada estrato es un estimador insesgado de la cuasivarianza poblacional (SR) y de la varianza poblacional (CR).

$$\text{Cuasiv. muestral} \rightarrow \hat{S}_h^2 = \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (X_{hi} - \bar{X}_h)^2 = \frac{n_h}{n_h - 1} \sigma_h^2$$

$$\text{Cuasiv. poblacional} \rightarrow S_h^2 = \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h)^2 = \frac{N_h}{N_h - 1} \sigma_h^2$$

En el caso de las proporciones, características cualitativas.

$\hat{S}_h^2 = \frac{n_h}{n_h - 1} \hat{P}_h \hat{Q}_h$ es un estimador insesgado de la cuasivarianza poblacional $S_h^2 = \frac{N_h}{N_h - 1} P_h Q_h$ (SR) y de la varianza poblacional $\sigma_h^2 = P_h Q_h$ (CR).

Para obtener las expresiones de las estimaciones de las varianzas basta con substituir cada parámetro poblacional por su estimador insesgado en la expresión de la varianza del estimador.

SIN reposición:

$$\hat{V}(\hat{X}_{st}) = \sum_{h=1}^L N_h^2 (1-f_h) \cdot \frac{\hat{S}_h^2}{n_h}$$

$$\hat{V}(\hat{\bar{X}}_{st}) = \sum_{h=1}^L W_h^2 (1-f_h) \cdot \frac{\hat{S}_h^2}{n_h}$$

$$\hat{V}(\hat{A}_{st}) = \sum_{h=1}^L N_h^2 (1-f_h) \cdot \frac{\hat{P}_h \hat{Q}_h}{n_h - 1}$$

$$\hat{V}(\hat{P}_{st}) = \sum_{h=1}^L W_h^2 (1-f_h) \cdot \frac{\hat{P}_h \hat{Q}_h}{n_h - 1}$$

CON reposición:

$$\hat{V}(\hat{X}_{st}) = \sum_{h=1}^L N_h^2 \cdot \frac{\hat{S}_h^2}{n_h}$$

$$\hat{V}(\hat{\bar{X}}_{st}) = \sum_{h=1}^L W_h^2 \cdot \frac{\hat{S}_h^2}{n_h}$$

$$\hat{V}(\hat{A}_{st}) = \sum_{h=1}^L N_h^2 \cdot \frac{\hat{P}_h \hat{Q}_h}{n_h - 1}$$

$$\hat{V}(\hat{P}_{st}) = \sum_{h=1}^L W_h^2 \cdot \frac{\hat{P}_h \hat{Q}_h}{n_h - 1}$$

3- AFIJACIÓN de la MUESTRA con UNA CARACTERÍSTICA

Se llama afijación de la muestra al reparto del tamaño muestral entre todos los estratos.

Consiste en determinar el valor de n_h , $h=1, \dots, L$, de modo que $n_1 + \dots + n_L = n$.

Los criterios de afijación más importantes son:

1- Afijación UNIFORME o igual

Consiste en asignar el mismo número de unidades muestrales a cada estrato:

$$n_h = \frac{n}{L} = k \quad \forall h=1, \dots, L \quad \rightarrow f_h = \frac{k}{N_h}$$

aumentando o disminuyendo en unidades si n no es múltiplo de L .

Este tipo de afijación ^{al} ~~no~~ ^{brings} importancia a todos los estratos, favorece a los estratos pequeños y perjudica a los grandes en cuanto a precisión. Sólo es conveniente en estratos de tamaño similar.

No interviene en absoluto el hecho de que el muestreo sea con o sin reposición.

2- Afijación PROPORCIONAL

Consiste en asignar a cada estrato un número de unidades muestrales proporcional a su tamaño

$$n_h = k N_h \Rightarrow n = \sum_h n_h = k \sum_h N_h = k N \Leftrightarrow k = \frac{n}{N}$$

$$n_h = \frac{n}{N} \cdot N_h$$

$$f_h = \frac{n_h}{N_h} = \frac{n}{N} = f$$

Directamente, se comprueba:

$$f_h = k = f$$

$$W_h = \frac{N_h}{N} = \frac{n_h/k}{n/k} = \frac{n_h}{n}$$

En el caso de afijación proporcional, las expresiones de los estimadores del total poblacional y de la media poblacional pueden:

$$\hat{X}_{st} = \sum_{h=1}^L N_h \bar{x}_h = \sum_{h=1}^L \left(\frac{n_h}{K} \right) x_h = \frac{1}{K} \sum_{h=1}^L n_h \cdot \frac{x_h}{n_h} = \frac{x}{K} = \frac{\text{Total muestral}}{\text{Fracción muestreo}}$$

$$\hat{\bar{X}}_{st} = \sum_{h=1}^L w_h \bar{x}_h = \sum_{h=1}^L \frac{n_h}{n} \bar{x}_h = \frac{1}{n} \sum_{h=1}^L n_h \cdot \frac{x_h}{n_h} = \frac{x}{n} = \frac{\text{Total muestral}}{\text{tamaño muestral}}$$

- Las fracciones de muestreo en los estratos son iguales y coinciden con la fracción global de muestreo, igual a la de proporcionalidad.
 - Los coef. de ponderación w_h se obtienen exclusivamente a partir de la muestra. $w_h = \frac{N_h}{N} = \frac{n_h/K}{n/K} = \frac{n_h}{n}$
 - Todas las unidades tienen la misma probabilidad de figurar en la muestra \rightarrow muestras autoponderadas.
 - $V(\bar{x}_{st}) = \frac{1-f}{n} \sum w_h s_h^2$ (??)
3. Afijación de MÍNIMA VARIANZA (o de Neyman)

Consiste en determinar los tamaños muestrales de cada estrato para un tamaño muestral fijo n que minimicen la varianza del estimador.

En todos los casos, se resuelve aplicando el mt. de los multiplicadores de Lagrange al problema de optimización con restricciones:

$$\begin{cases} \min_{n_h} V(\hat{\theta}_{st}) \\ \text{s.a.} \sum_{h=1}^L n_h = n \end{cases}$$

Al ser las expresiones de la varianza del estimador distintas para cada caso hay que resolver por separado, distinguiendo los casos SIN y CON reposición.

SIN reposición, para la media y la proporción

$$\begin{aligned}
 V(\hat{\bar{x}}_{st}) &= V(\bar{x}_{st}) = \sum_{h=1}^L W_h^2 (1-f_h) \cdot \frac{S_h^2}{n_h} = \sum_{h=1}^L W_h^2 \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{S_h^2}{n_h} = \\
 &= \sum_{h=1}^L W_h^2 \cdot \frac{S_h^2}{n_h} - \sum_{h=1}^L W_h^2 \cdot \frac{n_h}{N_h} \cdot \frac{S_h^2}{n_h} = \\
 &= \sum_{h=1}^L W_h^2 \cdot \frac{S_h^2}{n_h} - \sum_{h=1}^L W_h \cdot \frac{N_h}{N} \cdot \frac{S_h^2}{N} \\
 &\quad \text{no depende de } n_h
 \end{aligned}$$

$$\left. \begin{array}{l} \text{MIN}_{n_h} V(\bar{x}_{st}) \\ \text{s.a. } \sum_{h=1}^L n_h = n \end{array} \right\} \sim \left. \begin{array}{l} \text{MIN}_{n_h} \sum_{h=1}^L W_h^2 \cdot \frac{S_h^2}{n_h} \\ \text{s.a. } \sum_{h=1}^L n_h = n \end{array} \right\}$$

1. Lagrangiano: $\phi(n_h, \lambda) = \sum_{h=1}^L W_h^2 \cdot \frac{S_h^2}{n_h} + \lambda \left(\sum_{h=1}^L n_h - n \right)$

$$\left. \begin{array}{l} 2. \frac{\partial \phi}{\partial n_h} = -W_h^2 \cdot \frac{S_h^2}{n_h^2} + \lambda = 0 \\ \text{para cada } h=1, \dots, L \\ \frac{\partial \phi}{\partial \lambda} = \sum n_h - n = 0 \end{array} \right\} \Rightarrow \lambda = W_h^2 \cdot \frac{S_h^2}{n_h^2} = \left(\frac{N_h}{N} \cdot \frac{S_h}{n_h} \right)^2$$

$$\sqrt{\lambda} = \frac{N_h}{N} \cdot \frac{S_h}{n_h} \Rightarrow N\sqrt{\lambda} = \frac{N_h S_h}{n_h} \Rightarrow \frac{n_h}{N_h S_h} = \frac{1}{N\sqrt{\lambda}} = \text{cte}$$

Desarrollando la igualdad para todo h y aplicando las propiedades de las proporciones:

$$\frac{n_h}{N_h S_h} = \frac{\sum n_h}{\sum N_h S_h} = \frac{n}{\sum N_h S_h} \Rightarrow$$

$$n_h = n \cdot \frac{N_h S_h}{\sum_{h=1}^L N_h S_h} = n \cdot \frac{W_h S_h}{\sum W_h S_h}$$

- En el supuesto de que todos los estratos tengan la misma dispersión, $S_h = S$, la asignación mínima varianta coincide con la asignación proporcional.

$$W_h = \frac{N_h}{N}$$

- La afijación óptima resulta útil si hay grandes diferencias de variabilidad entre los estratos. En caso contrario, por su sencillez, se recomienda la afijación proporcional.

- El valor de la varianza mínima es:

$$V(\bar{x}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L W_h S_h \right)^2 - \frac{1}{N} \sum_{h=1}^L W_h S_h^2$$

- Si se quiere la expresión de la varianza y la afijación para el estimador de la proporción basta sustituir en las fórmulas anteriores S_h^2 por $\frac{N_h(P_h Q_h)}{N_h - 1}$.

- En el caso del total poblacional y del total de clase, cambia la función objetivo, y aplicando multiplicadores de Lagrange se llega a la misma afijación que en el caso de la media, porque la función objetivo se diferencia en una de $(N_h = W_h \cdot N)$.

- El valor de la varianza mínima para el estimador del total es

$$V(\hat{X}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L N_h S_h \right)^2 - \frac{1}{N} \sum_{h=1}^L N_h S_h^2$$

CON reposición, para la media y la proporción

$$\left. \begin{array}{l} \text{MIN } V(\bar{x}_{st}) \\ \text{s.a. } \sum n_h = n \end{array} \right\} , \text{ donde } V(\bar{x}_{st}) = \sum_{h=1}^L W_h^2 \cdot \frac{\sigma_h^2}{n_h}$$

1- Lagrangiano: $\phi(n_h, \lambda) = \sum_{h=1}^L W_h^2 \cdot \frac{\sigma_h^2}{n_h} + \lambda (\sum n_h - n)$

2- Derivar e igualar a cero:

$$\left. \begin{array}{l} \frac{\partial \phi}{\partial n_h} = -W_h^2 \cdot \frac{\sigma_h^2}{n_h^2} + \lambda = 0 \\ \frac{\partial \phi}{\partial \lambda} = \sum n_h - n = 0 \end{array} \right\} \begin{array}{l} \lambda = W_h^2 \cdot \frac{\sigma_h^2}{n_h^2} = \frac{N_h^2}{N^2} \cdot \frac{\sigma_h^2}{n_h^2} \\ \downarrow \\ \sqrt{\lambda} = W_h \cdot \frac{\sigma_h}{n_h} = \text{cte} \end{array}$$

Aplicando las propiedades de las proporciones:

$$\frac{W_h \cdot \sigma_h}{n_h} = \frac{\sum W_h \sigma_h}{\sum n_h} = \frac{\sum W_h \sigma_h}{n} \Rightarrow n_h = n \cdot \frac{W_h \sigma_h}{\sum_{h=1}^L W_h \sigma_h} =$$

$$\boxed{n_h = n \cdot \frac{N_h \sigma_h}{\sum_{h=1}^L N_h \sigma_h}} = n \cdot \frac{N_h \sigma_h}{\sum_{h=1}^L N_h \sigma_h}$$

- Los valores de n_h son proporcionales a los productos $N_h \cdot \sigma_h$. En el supuesto de que las dispersiones de los estratos coincidan, la afijación de mínima varianza coincide con la afijación proporcional.
- El valor de la varianza mínima es:

$$V(\bar{x}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L W_h \sigma_h \right)^2$$
- Si se fuere la afijación de mínima varianza y la expresión de varianza mínima para el estimador de la proporción basta con sustituir σ_h^2 por $P_h Q_h$.
- En el caso del total poblacional, el tamaño muestral de cada estrato coincide con la media.

El valor de la varianza mínima

4- Afijación ÓPTIMA con costes

Consiste en determinar los tamaños muestrales para cada estrato, n_h , $h=1 \dots L$, de forma que para un coste fijo C la varianza de los estimadores sea mínima.

El coste fijo C será la suma de los costes de cada unidad. $C = \sum_{h=1}^L c_h n_h$

Vuelve a ser un problema de optimización condicionado que se resuelve con los multiplicadores de Lagrange.

SIN reemplazamiento, para la media y la proporción

$$\left. \begin{array}{l} \text{MIN } V(\bar{x}_{st}) \\ \text{s.a. } \sum_{h=1}^L c_h n_h = C \end{array} \right\} \quad \text{donde } V(\bar{x}_{st}) = \sum_{h=1}^L w_h^2 \cdot \frac{S_h^2}{n_h} - \underbrace{\sum_{h=1}^L w_h \cdot \frac{S_h^2}{N}}_{\text{no depende}}$$

por lo que $\text{MIN } V(\bar{x}_{st}) \sim \text{MIN} \left(\sum_{h=1}^L w_h^2 \frac{S_h^2}{n_h} \right)$

1 - Lagrangiano: $\phi(n_h, \lambda) = \sum_{h=1}^L w_h^2 \cdot \frac{S_h^2}{n_h} + \lambda \left(\sum_{h=1}^L c_h n_h - C \right)$

2 - Derivando e igualando a 0:

$$\left. \begin{array}{l} \frac{\partial \phi}{\partial n_h} = -w_h^2 \cdot \frac{S_h^2}{n_h^2} + \lambda c_h = 0 \\ \frac{\partial \phi}{\partial \lambda} = \sum c_h n_h - C = 0 \end{array} \right\} \quad \lambda = \frac{w_h^2 \cdot \frac{S_h^2}{n_h^2}}{c_h} = \frac{1}{c_h} \cdot \frac{N_h^2}{N^2} \cdot \frac{S_h^2}{n_h^2}$$

$$\frac{\partial \phi}{\partial \lambda} = \sum c_h n_h - C = 0$$

$$\sqrt{\lambda} = \frac{N_h}{N} \cdot \frac{S_h}{\sqrt{c_h} \cdot n_h}$$

$$N\sqrt{\lambda} = \frac{N_h \cdot S_h / \sqrt{c_h}}{n_h} \Rightarrow$$

$$\Rightarrow \frac{n_h}{N_h S_h / \sqrt{c_h}} = \frac{1}{N\sqrt{\lambda}} = \text{cte}$$

Desarrollando la igualdad para todo h y aplicando las propiedades de las proporciones:

$$\frac{n_h}{N_h S_h / \sqrt{c_h}} = \frac{\sum n_h}{\sum N_h S_h / \sqrt{c_h}} = \frac{n}{\sum N_h S_h / \sqrt{c_h}} \Rightarrow$$

$$\Rightarrow n_h = n \cdot \frac{N_h S_h / \sqrt{c_h}}{\sum_{h=1}^L N_h S_h / \sqrt{c_h}} = n \cdot \frac{w_h S_h / \sqrt{c_h}}{\sum_{h=1}^L w_h S_h / \sqrt{c_h}}$$

- Los valores de n_h son proporcionales a los productos $N_h S_h / \sqrt{c_h}$, y en el supuesto de que el coste C es

$$- V_{\text{func}}(\bar{x}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L w_h S_h / \sqrt{c_h} \right) \left(\sum_{h=1}^L w_h S_h / \sqrt{c_h} \right) - \frac{1}{N} \sum_{h=1}^L w_h \cdot S_h^2$$

sea constante en todos los estratos, $c_h = k \forall h$, entonces la afijación óptima coincide con la de mínima varianza. Si además, las dispersiones de los estratos son iguales, la afijación óptima coincidirá con la afijación proporcional.

~~El valor de la varianza mínima en este caso es:~~

- Si se quiere la afijación óptima para el estimador de la proporción basta con sustituir S_h^2 por $\frac{N_h}{N_h - 1} \cdot P_h Q_h$

- En el caso del total poblacional, el tamaño muestral de cada estrato coincide con el obtenido para la media, al diferenciarse $V(\bar{X}_{st})$ de $V(\hat{X}_{st})$ en uso de la proporcionalidad.

~~En~~

con reemplazamiento, para la media y la proporción

$$\begin{aligned} \text{MIN } V(\bar{X}_{st}) \\ \text{s.a. } \sum c_h n_h = C \end{aligned} \quad \left\{ \begin{array}{l} \text{donde } V(\bar{X}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L N_h Q_h \right) \end{array} \right.$$

Aplicando multiplicadores de Lagrange se llega a:

$$n_h = n \cdot \frac{N_h Q_h / \sqrt{C_h}}{\sum_{h=1}^L N_h Q_h / \sqrt{C_h}}$$

- Se observa que n_h es proporcional a $N_h Q_h / \sqrt{C_h}$. Por tanto de la afijación óptima coincide con la mínima. Si las varianzas coinciden además, ~~por ser~~ igual que la prop.

- El valor de la ~~varianza~~ ^{varianza} es $\frac{1}{n} \left(\sum N_h Q_h \sqrt{C_h} \right) \left(\sum N_h Q_h / \sqrt{C_h} \right)$.

- Que el estimador del total, N_h coincide.

4- CONSTRUCCIÓN y NÚMERO de ESTRATOS

→ No está en César Pérez
→ CONJUNTO

No es fácil dar reglas fijas con respecto al número de estratos y sus límites, pues dependen en cada caso de

- la estructura de la población
- la información disponible
- los objetivos concretos en cada caso.

Para que el muestreo estratificado dé buenos resultados y sus estimadores sean más precisos, los estratos deben estar constituidos por unidades lo más homogéneas posibles, y tienen que ser lo más heterogéneos posibles entre sí.

En general, la precisión aumenta cuanto mayor sea el número de estratos bien constituidos, pero se complican los cálculos y la presencia de cada estrato ~~disminuye~~ en la muestra disminuye.

En todo caso, los estratos deben constituir una partición de la población: cada unidad poblacional debe pertenecer a uno y sólo uno de los estratos, y $\sum_{h=1}^L N_h = N$.

~~Deben de~~

A veces el orden en que aparecen las unidades en el marco implica una estratificación. Otras veces, la necesidad de efectuar estimaciones pequeñas. Generalmente la estratificación se hace por razones de eficiencia (mayor precisión y menor coste).

CONSTRUCCIÓN de los estratos

Cuando se estudia una sola característica de la población, la mejor ~~a elegir es~~ para estratificar es la distribución de frecuencias, pero suele ser descomodida. La siguiente mejor opción es la distrib. de frecuencias de otra variable altamente correlacionada con X .

Conocido el número de estratos L , el mínimo y el máximo valor de la característica en la población, se trata de encontrar los límites intermedios entre los estratos que minimicen la varianza del estimador.

En el caso concreto de la media muestral obtenida SIN reposición, su varianza es:

$$V(\bar{X}_{st}) = \sum_{h=1}^L W_h^2 (1-f_h) \frac{S_h^2}{n_h}$$

~~Para~~ si la muestra ha sido obtenida con afijación óptima de Neyman, suponiendo la fracción de muestreo en cada estrato despreciable:

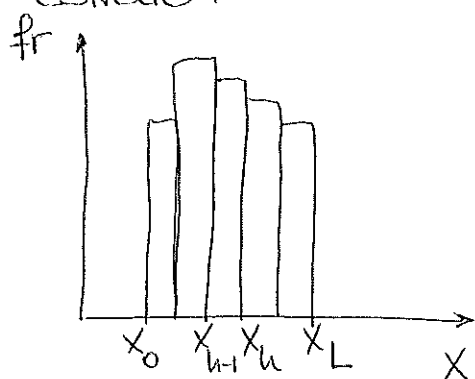
$$\left. \begin{aligned} n_h &= n \frac{W_h S_h}{\sum W_h S_h} \\ f_h &= \frac{n_h}{N_h} \approx 0 \end{aligned} \right\} \rightarrow V(\bar{X}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L W_h S_h \right)^2$$

por lo que minimizar la varianza del estimador equivale a minimizar $\sum_{h=1}^L W_h S_h$.

Dalenius y Hodges (1959) idearon un método rápido y aproximado basado en la cantidad

$$Z(X) = \int_{x_0}^X \sqrt{f(t)} dt \quad , \text{ donde } f(x) \equiv \text{función de frecuencias de } x$$

Según estos autores, si los estratos son muy numerosos y estrechos, la función de frecuencias debería ser aproximadamente constante (uniforme) dentro de cada estrato:



$$W_h = \int_{x_{h-1}}^{x_h} f_r(t) dt \simeq f_{r_h} (x_h - x_{h-1})$$

$$S_h \simeq \frac{1}{\sqrt{12}} (x_h - x_{h-1}) \quad , \quad \text{ap } N_{h-1} \simeq N_h$$

$$Z(x_h) - Z(x_{h-1}) = \int_{x_{h-1}}^{x_h} f_r(t) dt \simeq \sqrt{f_{r_h}} (x_h - x_{h-1})$$

Sustituyendo estas aproximaciones en la aproximación a minimizar resulta:

$$\sum_{h=1}^L W_h S_h \simeq \sum_{h=1}^L f_{r_h} \frac{(x_h - x_{h-1})^2}{\sqrt{12}} \simeq \sum_{h=1}^L \frac{(Z(x_h) - Z(x_{h-1}))^2}{\sqrt{12}}$$

que se minimiza al hacer $Z(x_h) - Z(x_{h-1})$ constante, por ser $(Z(x_L) - Z(x_0))$ una cantidad fija.

Así, la regla es elegir los x_h de modo que formen intervalos iguales en la escala $Z(x)\sqrt{f_r(x)}$ acumulada.

Es importante puntualizar que resulta poco realista basar la estratificación en los valores de la cart. X .

En la práctica, se utiliza alguna variable altamente correlacionada con X como variable estratificador a partir de la regresión lineal de X sobre ella.

NÚMERO de ESTRATOS

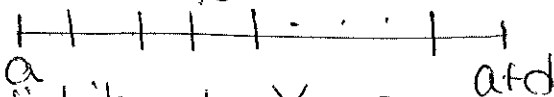
Si los estratos se constuyen con los valores de X , el caso más sencillo es que X siga una distrib. Uniforme en el intervalo $(a, a+d)$ de modo que:

$$\sigma_x^2 = \frac{d^2}{12} \Rightarrow S_x^2 = \frac{d^2}{12} \quad \text{c.p. } N-1 \approx N \text{ (N grande)}$$

Con una m.a.s. sin estratificar, la variancia del estimador media muestral sería (SIN reposición):

$$V(\bar{x}) = (1-f) \frac{S_x^2}{n} = \frac{d^2}{12n} \quad , \text{ sp } N \text{ grande} \Rightarrow 1-f \approx 1$$

Si se divide la población en L estratos de igual tamaño, la variancia dentro de cada uno de ellos será: $S_{x_h}^2 = \frac{d^2}{12L^2}$



pues dentro de cada estrato la distrib. de X es uniforme en un intervalo de amplitud $\frac{d}{L}$.

Por tanto,

$$\left. \begin{aligned} W_h &= \frac{d/L}{d} = \frac{1}{L} \\ n_h &= \frac{n}{L} \end{aligned} \right\} \Rightarrow V(\bar{x}_{st}) = \sum_{h=1}^L W_h^2 \frac{S_{x_h}^2}{n_h} = \sum_{h=1}^L \frac{1}{L^2} \cdot \frac{d^2}{12L^2} \cdot \frac{L}{n} = \frac{d^2}{12L^3 n} = \frac{V(\bar{x})}{L^3} \leftarrow \text{CTX } (2)$$

Entonces la variancia de \bar{x}_{st} decrece al aumentar el n.º de estratos \Rightarrow mejora la precisión.

Cuando se utiliza una var. estratificada, el aumento del n.º de estratos sólo es beneficioso si las var. están altamente correlacionadas.

En cualquier caso, hay que tener en cuenta que el aumento del número de estratos implica un aumento del coste, salvo que se reduzca el tamaño de la muestra, lo que no siempre compensa.

5. PRINCIPIOS BÁSICOS de la ESTRATIFICACIÓN

Si la estratificación es adecuada, debería obtenerse una varianza más pequeña para el estimador de μ que se obtendría con un muestreo aleatorio simple.

Para ver qué se entiende por "estratificación adecuada", se acude a la descomposición de la Varianza Total de la población en la suma de la variabilidad dentro de cada estrato y la variabilidad entre los estratos.

$$\begin{aligned}
 VT &= \sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X})^2 = \sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h + \bar{X}_h - \bar{X})^2 = \\
 &= \sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h)^2 + \sum_{h=1}^L \sum_{i=1}^{N_h} (\bar{X}_h - \bar{X})^2 + 2 \sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h)(\bar{X}_h - \bar{X}) \\
 &= \sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h)^2 + \sum_{h=1}^L N_h (\bar{X}_h - \bar{X})^2 + 2 \sum_{h=1}^L (\bar{X}_h - \bar{X}) \underbrace{\left(X_{h1} - \frac{N_h \bar{X}_h}{N_h} \right)}_0
 \end{aligned}$$

luego:

$$\sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X})^2 = \sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h)^2 + \sum_{h=1}^L N_h (\bar{X}_h - \bar{X})^2$$

$$\underbrace{VT}_{SCT} = \underbrace{\sum_{h=1}^L \sum_{i=1}^{N_h} (X_{hi} - \bar{X}_h)^2}_{VD} + \underbrace{\sum_{h=1}^L N_h (\bar{X}_h - \bar{X})^2}_{VE}$$

que también se puede expresar en términos de la varianza poblacional y de la varianza por estrato:

$$VT = VD + VE$$

$$(N-1)S_T^2 = (N-L)S_D^2 + (L-1)S_E^2$$

Partiendo de la expresión de la Varianza de la media muestral en muestreo aleatorio simple sin reposición, y substituyendo la varianza poblacional por la descomposición antes obtenida se llega a una expresión de donde se deduce el Principio Fundamental de la Estratificación.

$$V(\bar{X}) = (1-f) \frac{S^2}{n} = \frac{1-f}{n} \left[\underbrace{\frac{1}{N-1} \sum_h \sum_i (X_{hi} - \bar{X}_h)^2}_{\downarrow} + \frac{1}{N-1} \sum_h N_h (\bar{X}_h - \bar{X})^2 \right] =$$

$$\frac{1}{N-1} \cdot \sum_{h=1}^L (N_h - 1) S_h^2 = \underbrace{\frac{1}{N-1} \sum_h N_h S_h^2}_{\downarrow \frac{N}{N-1} \sum_h \frac{N_h}{N} S_h^2} - \frac{1}{N-1} \sum_{h=1}^L S_h^2$$

$$= \frac{N-1+1}{N-1} \sum_h \frac{N_h}{N} S_h^2 - \frac{1}{N-1} \sum_h S_h^2 = \sum_h \frac{N_h}{N} S_h^2 + \frac{1}{N-1} \sum_h \frac{N_h}{N} S_h^2 - \frac{1}{N-1} \sum_h S_h^2$$

$$= \sum_h \frac{N_h}{N} S_h^2 + \frac{1}{N-1} \sum_h \underbrace{\left(\frac{N_h}{N} - 1 \right)}_{\frac{N_h - N}{N}} S_h^2 = \sum_{h=1}^L \frac{N_h}{N} S_h^2 + \frac{1}{N(N-1)} \sum_{h=1}^L (N - N_h) S_h^2$$

por lo que

$$V(\bar{X}) = \frac{1-f}{n} \left[\sum_{h=1}^L \underbrace{\left(\frac{N_h}{N} \right)}_{W_h} S_h^2 - \frac{1}{N(N-1)} \sum_{h=1}^L (N - N_h) S_h^2 + \frac{1}{N-1} \sum_{h=1}^L N_h (\bar{X}_h - \bar{X})^2 \right] =$$

Para simplificar esta expresión, utilizando afinación proporcional:

$$\left. \begin{aligned} W_h &= \frac{N_h}{N} \stackrel{\text{prop}}{=} \frac{n_h}{n} \\ f_h &= \frac{n_h}{N_h} = \frac{n}{N} \end{aligned} \right\} \rightarrow V_{\text{prop}}(\bar{X}_{\text{st}}) = \sum_{h=1}^L W_h \frac{S_h^2}{n} (1-f)$$

luego

$$V(\bar{X}) = V_{\text{prop}}(\bar{X}_{\text{st}}) + \frac{1-f}{n(N-1)} \left[\frac{1}{N-1} \sum_{h=1}^L N_h (\bar{X}_h - \bar{X})^2 - \sum_{h=1}^L (N - N_h) \frac{S_h^2}{N} \right]$$

$$V(\bar{X}) = V_{\text{prop}}(\bar{X}_{\text{st}}) + \frac{1-f}{n(N-1)} \left[\underbrace{\sum_{h=1}^L N_h (\bar{X}_h - \bar{X})^2}_{\text{GRANDE}} - \underbrace{\frac{1}{N} \sum_{h=1}^L (N - N_h) S_h^2}_{\text{PEQUEÑO}} \right]$$

\downarrow
 V_{mas}

\downarrow
 $V_{\text{estadist. af. prop}}$

\downarrow
cantidad que puede ser positiva o negativa, por lo que no siempre es mejor

de donde se deduce el Principio Fundamental de la Estratificación:

"La ganancia en precisión con la estratificación será tanto mayor cuanto mayor sea el término entre corchetes, es decir, los estratos deben constituirse de manera qd:

- las medias de los estratos difieran lo más posible
- Dentro de cada estrato exista una gran homogeneidad entre sus unidades.

6- GANANCIA de PRECISIÓN

En el supuesto de poblaciones grandes con estratos grandes ($N-1 \approx N$ y $N_h-1 \approx N_h$) podemos descomponer la variancia poblacional en

$$S^2 \approx \sum_h W_h S_h^2 + \sum_h W_h (\bar{X}_h - \bar{X})^2$$

por lo que:

$$(1-f) \frac{S^2}{n} = \underbrace{\frac{1-f}{n} \sum_h W_h S_h^2}_{V_{Af.Prop}(\bar{X}_{st})} + \underbrace{\frac{1-f}{n} \sum_h W_h (\bar{X}_h - \bar{X})^2}_{\geq 0}$$

$$V_{mas}(\bar{X}) = V_{Af.Prop}(\bar{X}_{st}) + \geq 0$$

Acabamos de demostrar que el muestreo estratificado con afijación proporcional es más preciso que el muestreo aleatorio simple (SR).

Serán iguales cuando las medias de los estratos sean iguales (pobl. homogénea), por lo que la ganancia en precisión será mayor cuanto más distintas entre sí sean las medias de los estratos.

Ahora vamos a comparar las precisiones de la afijación proporcional y la de mínima varianza:

$$V_{AfProp}(\bar{x}_{st}) - V_{MinVar}(\bar{x}_{st}) = \underbrace{\frac{1-f}{n} \sum_{h=1}^L W_h S_h^2}_{\frac{1}{n} - \frac{1}{N}} - \left[\frac{1}{n} \left(\sum W_h S_h \right)^2 - \frac{1}{N} \sum W_h S_h^2 \right]$$

$$\begin{aligned} &= \frac{1}{n} \left[\sum W_h S_h^2 - \underbrace{\left(\sum W_h S_h \right)^2}_{\bar{S}} \right] + \frac{1}{N} \left[\underbrace{\sum W_h S_h^2}_{0} - \sum W_h S_h^2 \right] = \\ &= \frac{1}{n} \sum_{h=1}^L W_h (S_h - \bar{S})^2 \geq 0 \quad \& \end{aligned}$$

luego el muestreo estratificado con afijación de mínima varianza es más preciso que el muestreo con afijación proporcional (se da la igualdad cuando las desviaciones típicas de todos los estratos son iguales).

la ganancia en precisión de muestreo estratificado con af. MV ~~será~~ respecto al m.e. af. prop. será tanto mayor cuanto más ~~distinto~~ heterogéneos entre sí sean los estratos,

En general : Af. Min. Var \gg Af. Proporc \gg MAS

- Los resultados son iguales para los estimadores de X, A, P , ya que las varianzas difieren en constantes de proporcionalidad.

- En el caso de muestreo con reemplazamiento se llega a la misma conclusión, a partir de la expresión de la varianza poblacional.