

SÔBRE UMA DISTRIBUIÇÃO
DISCRETA DE
PROBABILIDADES

Tese apresentada no concurso para
provimento efetivo da CADEIRA
VI = ESTATÍSTICA I = da
Faculdade de Filosofia, Ciências e
Letras da Universidade de São Paulo.

SÃO PAULO

≡ 1952 ≡

SÔBRE UMA DISTRIBUIÇÃO
DISCRETA DE PROBABILIDADES

Tese apresentada no concurso para
provimento efetivo da CADEIRA VI
- ESTATÍSTICA I - da Faculdade de
Filosofia, Ciências e Letras da U-
niversidade de São Paulo

INTRODUÇÃO

Sempre tivemos para nós que a inclusão da prova de tese em um concurso visa precìpuaamente permitir à Comissão Julgadora aquilatar a eventual originalidade científica do candidato.

Este ponto de vista constitui, no âmbito estatístico, para aquêles que têm o ensino como missão primordial, um sério óbice. Efetivamente, pela própria natureza e condições do seu trabalho, são êles conduzidos, quasi forçosamente, a perquirir o algo de novo exigido, no domínio puramente teórico. Ora, em um ambiente como o nosso - que é forçoso reconhecer ainda incipiente em questões estatísticas - o conseguir oportunidade para realizar qualquer coisa de inédito em matéria doutrinária é tarefa que demanda não só assiduidade e perseverança no estudo, como tempo de duração difficilmente previsível. Se isto não bastasse para arrefecer o entusiasmo de pesquisadores honestos em empreender a grande caminhada, teríamos ainda de lembrar o entrave decorrente de escassez bibliográfica que dificulta sobremodo a cada qual objetivar, de forma positiva, a convicção de originalidade que êle, em sua consciência, empresta, por vêzes, a certos resultados que consegue.

As considerações supra visam realçar a sa -

tisfação íntima que experimentamos quando, ao estudar certos problemas de probabilidade, imaginamos uma questão de amostragem de postos que nos conduzia a uma distribuição discreta que nos pareceu, desde logo, digna de um estudo acurado.

Começamos, então, a viver sincera e intensamente o assunto, dando de nós mesmos os melhores e mais inteligentes esforços. Como, a medida que prosseguiam as nossas indagações, fossem aparecendo certos resultados interessantes, seja de natureza matemática (por exemplo, o encontro de uma expressão geral da soma das potências de grau l dos n primeiros números naturais) seja de natureza propriamente estatística (v.g. a generalização da distribuição para o caso multidimensional e o que é ainda mais importante, a prova do teorema básico dos testes não paramétricos), persuadimo-nos que os frutos da tenacidade do nosso labor mereciam as honras de corporificar uma tese destinada a concorrer ao mais alto posto na Faculdade que, sem dúvida, é a cúpula da Universidade de São Paulo.

1 - O problema

Uma população finita contém n indivíduos ordenados de 1 a n segundo um certo critério. Toma-se uma amostra casual de k ($k \leq n$) indivíduos, sem reposição, cograduando-os segundo o posto por eles ocupado na população.

Pergunta-se : Qual é a probabilidade do indivíduo que na população ocupa o posto y vir a ocupar na amostra o posto h ($1 \leq h \leq k$) ?

Para que o indivíduo de posto y na população ocupe o posto h na amostra é necessário e suficiente que na amostra apareçam $h-1$ indivíduos com posto, na população, menor do que y , o indivíduo de posto y e $k-h$ indivíduos com posto, na população, maior do que y .

Ora, $h-1$ indivíduos com posto menor do que y podendo ser escolhidos de $\binom{y-1}{h-1}$ diferentes maneiras e $k-h$ com posto maior do que y de $\binom{n-y}{k-h}$ diferentes maneiras, então, o número de casos favoráveis é :

$$\binom{y-1}{h-1} \binom{n-y}{k-h}$$

quantidade que somente é diferente de zero para :

$$y = h, h+1, \dots, (n-k+h)$$

Assim sendo, e desde que o número de amostras cograduadas é $\binom{n}{k}$, a probabilidade procurada é

- que indicaremos com $f(y)$ - vale :

$$f(y) = \begin{cases} \frac{\binom{y-1}{k-1} \binom{n-y}{k-h}}{\binom{n}{k}} & \text{se } y = h, h+1, \dots, n-k+h \\ 0 & \text{em contrário} \end{cases} \quad (1)$$

A $f(y)$ depende dos parâmetros n , h e k , tomando-se este vocábulo no sentido matemático e não no de parâmetros populacionais. Apresenta ela certa analogia formal com a distribuição hipergeométrica, da qual difere, entretanto, pelo fato de que, enquanto nesta, a variável é a classe da combinação (denominador), naquela, a variável é o número de elementos (numerador).

A $f(y)$ dá, naturalmente, a distribuição da variável aleatória "pôsto populacional que ocupa na amostra o pôsto h " quando se tomam, sem reposição, de forma casual, amostras de tamanho k de uma população constituída de n pôstos.

Como a grandeza pôsto pode ser equiparada a uma variável definida no campo dos números naturais, a interpretação supra se traduz, dizendo que a $f(y)$ dá a distribuição amostral da variável aleatória "valor que ocupa o pôsto h " quando se tomam, sem reposição, amostras casuais de tamanho k de uma população constituída pelos n valores $1, 2 \dots n$.

Por último, salientemos que a nossa questão difere do problema clássico da determinação da distribuição amostral do valor que ocupa um certo pôsto, nisto que, en-

quanto aqui a amostragem é sem reposição, no problema clássico, a amostra é simples.

A primeira vista, poderá causar estranheza a natureza da nossa variável aleatória. Esta impressão, entretanto, se desvanecerá lembrando que no problema básico das provas não paramétricas, a variável é dada por $\int_{x_r}^{x_s} f(x) dx$ ou seja, o total de área da distribuição populacional compreendida entre os dois valores x_r e x_s que ocupam na graduatória amostral os postos r e s ($s > r$).

Como primeira importante consequência, notemos que a (1) definindo uma função de frequência, a soma dos valores de $f(y)$ estendida ao campo de variação de y vale 1. Assim sendo, a (1) dá uma demonstração estatística da fórmula matemática

$$\sum_{z=0}^{n-k} \binom{z+k-1}{k-1} \binom{n-z-k}{k-k} = \binom{n}{k}$$

obtida da $f(y)$ pela transformação $y-h = z$.

2 - Determinação dos momentos

É fácil depreender, por uma análise ainda que perfunctória da (1), que a determinação dos momentos de potência da variável y apresenta sérias dificuldades. Pensamos, desde logo, em remover este impecilho, calculando as referidas características por via dos momentos fato-

riais. Novamente, aqui, deparava-se uma dificuldade consistente em que $y^{[j]}$ não se prestava a simplificações convenientes com os fatoriais constantes da expressão de $E(y^{[j]})$. Notando que no denominador desta última aparece $(n-y-k+h)!$, pensamos, então, começar a realização da tarefa constante desta epígrafe, procurando o valor de

$$E\left[(n-y-k+h)^{[j]}\right]$$

o que se processou como segue:

$$\begin{aligned} \alpha'_{[j]} &= E\left[(n-y-k+h)^{[j]}\right] \\ &= \frac{1}{\binom{n}{k}} \sum_{y=k}^{n-k+h} \binom{y-1}{h-1} (n-y-k+h)^{[j]} \binom{n-y}{h-k} \\ &= \frac{1}{\binom{n}{k}} \sum_{y=k}^{n-k+h} \binom{y-1}{h-1} \frac{(n-y)!}{(n-y-k+h-j)!(h-k)!} \\ &= \frac{(k-h+j)^{[j]}}{\binom{n}{k}} \sum_{y=k}^{n-k+h} \binom{y-1}{h-1} \frac{(n-y)!}{(n-y-k+h-j)!(h-k+j)!} \\ &= \frac{(k-h+j)^{[j]}}{\binom{n}{k}} \sum_{y=k}^{n-k+h} \binom{y-1}{h-1} \binom{n-y}{k+j-h} \end{aligned}$$

$$\begin{aligned}
 &= \frac{(k-h+d)^{[j]}}{\binom{n}{k}} \sum_{y=h}^{n-(k+j)+k} \binom{y-1}{k-1} \binom{n-y}{k+j-k} \\
 &= \frac{(k-h+d)^{[j]}}{\binom{n}{k}} \binom{n}{k-j} \\
 &= \frac{(k-h+d)^{[j]} (n-k)^{[j]}}{(k+j)^{[j]}}
 \end{aligned}$$

Utilizando, agora, a relação entre momentos de potências e momentos fatoriais, em torno de mesma origem, dada por

$$\alpha'_t = \sum_{j=0}^t T_{tj} \alpha_{[j]}$$

podemos escrever para o momento de potência de ordem t , não centrado, da variável $n-y-k+h$:

$$\alpha'_t = \sum_{j=0}^t T_{tj} \frac{(k-h+d)^{[j]} (n-k)^{[j]}}{(k+j)^{[j]}} \quad (2)$$

onde os T_{tj} são números de Stirling de segunda espécie.

Da (2) segue:

$$E(-y)^t = \sum_{j=0}^t T_{tj} \frac{(n+j)^{[j]} (-k)^{[j]}}{(k+j)^{[j]}}$$

e, portanto, simbolizando por α_l o momento não centrado, de ordem l , da variável y :

$$\alpha_l = (-1)^l \sum_{j=0}^l T_{lj} \frac{(n+j)^{[l]} (-k)^{[j]}}{(k+j)^{[j]}} \quad (3)$$

Notando que os momentos centrados de ordem l da variável y , a menos de $(-1)^l$ coincidem com os da variável $n-y-k+h$, então, simbolizando-os por μ_l , de (2) segue também:

$$\begin{aligned} \mu_l &= E \left[(n-y-k+h) - \alpha_l \right]^l \\ &= \sum_{\lambda=0}^l \binom{l}{\lambda} (-1)^\lambda (\alpha_l)^{l-\lambda} \alpha_\lambda \\ &= \sum_{\lambda=0}^l \sum_{j=0}^{\lambda} \binom{l}{\lambda} (-1)^\lambda T_{\lambda j} \left[\frac{(k-k+l)(n-k)}{(k+1)} \right]^{l-\lambda} \frac{(k-k+j)^{[j]} (n-k)^{[j]}}{(k+j)^{[j]}} \quad (4) \end{aligned}$$

Como casos particulares merecedores de registro explícito, estabeleçamos:

$$\chi_1 = \frac{(n+1)k}{k+1} \quad (5)$$

$$\mu_2 = \sum_{\lambda=0}^2 \sum_{j=0}^{\lambda} \binom{2}{\lambda} (-1)^\lambda T_{\lambda j} \left[\frac{(k-h+1)(n-k)}{k+1} \right]^{2-\lambda} \frac{(k-h+j)^{[j]} (n-k)^{[j]}}{(k+j)^{[j]}}$$

$$\begin{aligned}
 &= \left[\frac{(k-n+1)(n-k)}{k+1} \right]^2 - 2 \left[\frac{(k-n+1)(n-k)}{k+1} \right] + \\
 &+ \frac{(k-n+1)(n-k)}{k+1} + \frac{(k-n+2)(k-n+1)(n-k)(n-k-1)}{(k+2)(k+1)} \\
 &= \frac{k(n+1)(k-k+1)(n-k)}{(k+1)^2(k+2)} \quad (5)
 \end{aligned}$$

$$\begin{aligned}
 \mu_3 &= \sum_{i=0}^3 \sum_{j=0}^i \binom{3}{i} (-1)^j T_{ij} \left[\frac{(k-n+1)(n-k)}{k+1} \right]^{3-i} \frac{(n-k)^{[i]} (n-k)^{[j]}}{(k+j)^{[i]}} \\
 &= \left[\frac{(k-n+1)(n-k)}{k+1} \right]^3 - 3 \left[\frac{(k-n+1)(n-k)}{k+1} \right] + 3 \left[\frac{(k-n+1)(n-k)}{k+1} \right]^2 \\
 &+ 3 \frac{(k-n+2)(k-n+1)^2 (n-k)^2 (n-k-1)}{(k+2)(k+1)^2} - \frac{(k-n+1)(n-k)}{k+1} \\
 &- 3 \frac{(k-n+2)(k-n+1)(n-k)(n-k-1)}{(k+2)(k+1)} \\
 &- \frac{(k-n+3)(k-n+2)(k-n+1)(n-k)(n-k-1)(n-k-2)}{(k+3)(k+2)(k+1)}
 \end{aligned}$$

$$= -\frac{(k-k+1)(n-k)}{(k+1)^3(k+2)(k+3)} \left\{ (k-k+1)^2 \left[2(n-k)^2(k+2)(k+3) - \right. \right. \\ \left. \left. - 3(n-k)(n-k-1)(k+1)(k+3) + (n-k-1)(n-k-2)(k+1)^2 \right] + \right. \\ \left. + 3(k-k+1)(k+1) \left[-(n-k)(k+2)(k+3) - (n-k)(n-k-1)(k+3) + \right. \right. \\ \left. \left. + (n-k-1)(k+1)(k+3) + (n-k-1)(n-k-2)(k+1) \right] + \right. \\ \left. + (k+1)^2 \left[(k+2)(k+3) + 3(n-k-1)(k+3) + 2(n-k-1)(n-k-2) \right] \right\}$$

$$= -\frac{(k-k+1)(n-k)}{(k+1)^3(k+2)(k+3)} \left\{ (k-k+1)^2 \left[2(n-k)(k+3) \left\{ (n-k)(k+2) - (n-k-1)(k+1) \right\} + \right. \right. \\ \left. \left. + (n-k-1)(k+1) \left\{ -(n-k)(k+3) + (n-k-2)(k+1) \right\} \right] + \right. \\ \left. + 3(k-k+1)(k+1) \left[-(n-k)(k+3) \left\{ (k+2) + (n-k-1) \right\} + \right. \right. \\ \left. \left. + (n-k-1)(k+1) \left\{ (k+3) + (n-k-2) \right\} \right] + \right. \\ \left. + (k+1)^2 \left[(k+3) \left\{ (k+2) + (n-k-1) \right\} + 2(n-k-1) \left\{ (k+3) + (n-k-2) \right\} \right] \right\}$$

$$= -\frac{(k-k+1)(n-k)(n+1)}{(k+1)^3(k+2)(k+3)} \left\{ (k-k+1)^2 \left[2(n-k)(k+3) - 2(n-k-1)(k+1) \right] + \right. \\ \left. + 3(k-k+1)(k+1) \left[-(n-k)(k+3) + (n-k-1)(k+1) \right] + \right. \\ \left. + (k+1)^2 \left[(k+3) + 2(n-k-1) \right] \right\} =$$

$$= -\frac{(k-k+1)(n-k)(n+1)(k-k+1)}{(k+1)^3(k+2)(k+3)} \left[2(k-k+1)^2 - 3(k-k+1)(k+1) + (k+1)^2 \right]$$

donde, finalmente :

$$u_3 = - \frac{(k-n+1)(n-k)(n+1)k(2n-k+1)(2k-k-1)}{(k+1)^3(k+2)(k+3)} \quad (7)$$

Para o momento quarto centrado, tem-se :

$$\begin{aligned} \mu_4 &= \sum_{j=0}^4 \sum_{i=0}^j \binom{4}{i} (-1)^i T_{i,j} \left[\frac{(k-k+1)(n-k)}{(k+1)} \right]^{4-i} \frac{(k-k+j)^{[j]} (n-k)^{[i]}}{(k+j)^{[j]}} = \\ &= \frac{(k-k+1)^4 (n-k)^4}{(k+1)^4} - \frac{4(k-k+1)^4 (n-k)^4}{(k+1)^4} + \\ &+ 6 \frac{(k-k+1)^2 (n-k)^2}{(k+1)^2} \left[\frac{(k-k+1)(n-k)}{k+1} + \frac{(k-k+2)(k-k+1)(n-k)(n-k-1)}{(k+1)(k+2)} \right] - \\ &- 4 \frac{(k-k+1)(n-k)}{k+1} \left[\frac{(k-k+1)(n-k)}{k+1} + 3 \frac{(k-k+2)(k-k+1)(n-k)(n-k-1)}{(k+1)(k+2)} + \right. \\ &\quad \left. + \frac{(k-k+3)(k-k+2)(k-k+1)(n-k)(n-k-1)(n-k-2)}{(k+1)(k+2)(k+3)} \right] + \\ &+ \frac{(k-k+1)(n-k)}{k+1} + 7 \frac{(k-k+2)(k-k+1)(n-k)(n-k-1)}{(k+1)(k+2)} + \\ &+ 6 \frac{(k-k+3)(k-k+2)(k-k+1)(n-k)(n-k-1)(n-k-2)}{(k+1)(k+2)(k+3)} + \\ &+ \frac{(k-k+4)(k-k+3)(k-k+2)(k-k+1)(n-k)(n-k-1)(n-k-2)(n-k-3)}{(k+1)(k+2)(k+3)(k+4)} \end{aligned}$$

$$\begin{aligned}
 u_4 = & \frac{(k-k+1)(n-k)}{(k+1)^4(k+2)(k+3)(k+4)} \left\{ (k-k+1)^2 \left[-3(n-k)^3(k+2)(k+3)(k+4) + \right. \right. \\
 & + 6(n-k)^2(n-k-1)(k+1)(k+3)(k+4) - \\
 & - 4(n-k)(n-k-1)(n-k-2)(k+1)^2(k+4) + \\
 & \left. \left. + (n-k-1)(n-k-2)(n-k-3)(k+1)^3 \right] + \right. \\
 & + 6(k-k+1)^2(k+1) \left[(n-k)^2(k+2)(k+3)(k+4) + \right. \\
 & + (n-k)^3(n-k-1)(k+3)(k+4) - \\
 & - 2(n-k)(n-k-1)(k+1)(k+3)(k+4) - \\
 & - 2(n-k)(n-k-1)(n-k-2)(k+1)(k+4) + \\
 & + (n-k-1)(n-k-2)(k+1)^2(k+4) + \\
 & \left. \left. + (n-k-1)(n-k-2)(n-k-3)(k+1)^3 \right] + \right. \\
 & + (k-k+1)(k+1)^2 \left[-4(n-k)(k+2)(k+3)(k+4) - \right. \\
 & - 12(n-k)(n-k-1)(k+3)(k+4) - \\
 & - 8(n-k)(n-k-1)(n-k-2)(k+4) + \\
 & + 7(n-k-1)(k+1)(k+3)(k+4) + \\
 & + 10(n-k-1)(n-k-2)(k+1)(k+4) + \\
 & \left. \left. + 11(n-k-1)(n-k-2)(n-k-3)(k+1) \right] + \right. \\
 & + (k-1)^3 \left[(k+2)(k+3)(k+4) + 7(n-k-1)(k+3)(k+4) + \right. \\
 & + 10(n-k-1)(n-k-2)(k+4) + \\
 & \left. \left. + 6(n-k-1)(n-k-2)(n-k-3) \right] \right\}
 \end{aligned}$$

$$\begin{aligned}
 \mu_k &= \frac{(k-k+1)(n-k)}{(k-1)^2(k+2)(k+3)(k+4)} \left\{ (k-k+1)^2 \left[3(n-k)^2(k+3)(k+4) - (k+2)(n-k) - (k-1)(n-k+2) \right] \right. \\
 &\quad + 2(k-k)(n-k+2)(k-1)(k+4) \left\{ (n-k)(k+3) - (n-k+2)(k+1) \right\} \\
 &\quad \left. - (n-k+1)(n-k+2)(k+1)^2 \left\{ (n-k)(k+4) - (n-k+3)(k+2) \right\} \right\} \\
 &\quad + 6(k-k+1)^2(k+1) \left\{ (n-k)^2(k+3)(k+4) \left\{ (k+2) + (n-k-1) \right\} \right. \\
 &\quad \left. - 2(n-k)(n-k+1)(k+1)(k+4) \left\{ (k+1) + (n-k-2) \right\} \right. \\
 &\quad \left. + (n-k+1)(n-k+2)(k+1)^2 \left\{ (k+1) + (n-k-3) \right\} \right\} \\
 &\quad + (k-k+1)(k+1)^2 \left[4(n-k)(k+3)(k+4) \left\{ (k+2) + (n-k-1) \right\} \right. \\
 &\quad \left. - 2(n-k)(n-k+1)(k+4) \left\{ (k+3) + (n-k-2) \right\} \right. \\
 &\quad \left. + 2(n-k+1)(k+1)(k+4) \left\{ (k+3) + (n-k-2) \right\} \right. \\
 &\quad \left. + 12(n-k+1)(n-k+2)(k+2) \left\{ (k+4) + (n-k-3) \right\} \right] \\
 &\quad + (k+1)^2 \left\{ (k+3)(k+4) \left\{ (k+2) + (n-k-1) \right\} \right. \\
 &\quad \left. + 6(n-k+1)(k+4) \left\{ (k+3) + (n-k-2) \right\} \right. \\
 &\quad \left. + 6(n-k+1)(n-k+2) \left\{ (k+4) + (n-k-3) \right\} \right\}
 \end{aligned}$$

$$\begin{aligned}
 u_1 &= \frac{(k+1)(n-k)(n+1)}{(k+1)^2(k+2)(k+3)(k+4)} \left\{ 3(k-k+1)^2 \left[(n-k)(k+4)(2n-k+1) - (n-k-1)(k+1)(3n-k+2) \right] \right. \\
 &\quad + 6(k-k+1)^2(k+1) \left[(n-k)(k+4)(2n-k+1) - (n-k-1)(k+1)(3n-k+2) \right] \\
 &\quad + (k-k+1)(k+1)^2 \left[4(n-k)(k+4)(2n-k+1) + (n-k-1)(k+1)(3n-k+2) \right] \\
 &\quad \left. + (k+1)^2 \left[(k+4)(2n-k+1) + 2(n-k-1)(3n-k+2) \right] \right\} \\
 &= \frac{(k-k+1)(n-k)(n+1) k_0}{(k+1)^2(k+2)(k+3)(k+4)} \left\{ -3k(k-k+1) \left[(n-k)(k+4)(2n-k+1) - (n-k-1)(k+1)(3n-k+2) \right] \right. \\
 &\quad \left. + (k+1)^2 \left[(2n-k)(2n-k+1) + 2(n-k)(n+2) \right] \right\} \quad (8)
 \end{aligned}$$

Dêstes resultados seguem, por seu turno:

$$\gamma_1^2 = \frac{u_3^2}{u_2^2} = \frac{(2k-k-1)^2 (2n-k+1)^2 (k+2)}{k_0 (n+1)(k-k+1)(n-k_0)(k+3)^2}$$

$$\begin{aligned}
 \gamma_2^2 &= \frac{u_4}{\lambda k_0^2} = \frac{k+2}{(k+3)(k+4)(n+1)k_0(k-k+1)(n-k)} \left\{ -3k(k-k+1) \left[(n-k)(k+4)(2n-k+1) - \right. \right. \\
 &\quad \left. \left. - (n-k-1)(k+1)(3n-k+2) \right] + \right. \\
 &\quad \left. + (k+1)^2 \left[(2n-k)(2n-k+1) + \right. \right. \\
 &\quad \left. \left. + 2(n-k)(n+2) \right] \right\}
 \end{aligned}$$

Os resultados obtidos mostram, pois, que se $n(h-1)$ for múltiplo de $k-1$, a distribuição admite dois máximos, em correspondência aos valores $\frac{n(h-1)}{k-1}$ e $\frac{n(h-1)}{k-1} + 1$; se, ao contrário, $n(h-1)$ não for múltiplo de $k-1$, a distribuição admite um único máximo em correspondência ao inteiro compreendido entre os dois limites das desigualdades supra.

b)- Pondo na (1) $y=h$ e considerando o resultado como função de h , tem-se :

$$\varphi(h) = \frac{\binom{n-h}{k-h}}{\binom{n}{h}} \quad h = 1, 2, \dots, \frac{n}{k}$$

que dá a probabilidade de um elemento populacional de posto $\leq k$ vir a ocupar na amostra o mesmo posto que ele ocupa na população. É fácil mostrar que $\varphi(h)$ é decrescente e conseqüentemente máxima no ponto $h = 1$. Com efeito, tem-se :

$$\binom{n-h}{k-h} = \binom{n-h-1}{k-h-1} + \binom{n-h-1}{k-h} > \binom{n-[h+1]}{k-[h+1]}$$

c) - Em particular, se $k = n$, isto é, se a amostra é a própria população, y admite o único valor h e a distribuição torna-se :

$$f(y) = \begin{cases} 1 & \text{se } y = h \\ 0 & \text{em contrário} \end{cases}$$

que nada mais é do que uma distribuição t.

O fato mais interessante a assinalar, no caso em aprêço, é que da (6) segue:

$$\mu_2 = 0$$

como deveria acontecer, posto que, como vimos, estamos diante de uma distribuição com toda a massa concentrada em um único ponto.

d) - Para $k = 1$ e conseqüentemente $h = 1$, o nosso problema reduz-se à determinação da probabilidade de se retirar um único elemento da população. A distribuição torna-se :

$$f(y) = \begin{cases} \frac{1}{n} & \text{se } y = 1, 2 \dots n \\ 0 & \text{em contrário.} \end{cases}$$

o que mostra que a nossa distribuição contém, como caso particular, a distribuição uniforme.

Assim sendo, pareceu-nos de importância aproveitar alguns dos resultados obtidos para estabelecer outras tantas propriedades desta distribuição que tamanho papel desempenha seja em estatística seja no cálculo de probabilidades.

Neste sentido, citemos, de início, que da (3) resulta :

$$n x_l = \sum_{y=1}^n y^l = (-1)^l \sum_{j=0}^l T_{lj} \frac{(n+j)^{[l+1]} (-1)^{[j]}}{(j+1)^{[j]}} \quad (10)$$

com o que se obtém, de forma deveras elegante e elementar, uma expressão geral para a soma das l -ésimas potências dos n primeiros números naturais.

Em particular, fazendo $l = 1, 2, 3$ e 4 , dá (10) seguem as fórmulas clássicas :

$$\sum_{y=1}^n y = - \sum_{j=0}^1 T_{1j} \frac{(n+j)^{[1+1]} (-1)^{[j]}}{(j+1)^{[j]}} = \frac{n(n+1)}{2}$$

$$\sum_{y=1}^n y^2 = \sum_{j=0}^2 T_{2j} \frac{(n+j)^{[2+1]} (-1)^{[j]}}{(j+1)^{[j]}} = \frac{n(n+1)(2n+1)}{6}$$

$$\sum_{y=1}^n y^3 = - \sum_{j=0}^3 T_{3j} \frac{(n+j)^{[3+1]} (-1)^{[j]}}{(j+1)^{[j]}} = \frac{n^2 (n+1)^2}{4}$$

$$\sum_{y=1}^n y^4 = \sum_{j=0}^4 T_{4j} \frac{(n+j)^{[4+1]} (-1)^{[j]}}{(j+1)^{[j]}} = \frac{n(n+1)(2n+1)(3n^2+3n-1)}{30}$$

Outro importante resultado diz respeito aos momentos centrados em torno da média \bar{y} dos n primeiros números naturais.

Tem-se, genericamente pela (4) :

$$\frac{1}{n} \sum_{y=1}^n (y - \bar{y})^k = \sum_{\lambda=0}^k \sum_{j=0}^{\lambda} \binom{\lambda}{\lambda} (-1)^{\lambda} T_{\lambda j} \left(\frac{n-1}{2}\right)^{\lambda-j} \frac{j^{\lfloor \lambda/2 \rfloor} (n-1)^{\lfloor \lambda/2 \rfloor}}{(j+1)^{\lfloor \lambda/2 \rfloor}}$$

e, em particular pelas (6), (7) e (8) :

$$\mu_2 = \frac{n^2 - 1}{12}$$

$$\mu_3 = 0$$

$$\mu_4 = \frac{(n^2 - 1)(3n^2 - 7)}{240}$$

e) - Sob a hipótese $h = 1$ e $k > 1$, tem-se para a distribuição do extremo superior da graduatória amostral :

$$f(y) = \begin{cases} \frac{k(n-k)^{\lfloor y-1 \rfloor}}{n^{\lfloor y \rfloor}} & \text{se } y = 1, 2, \dots, (n-k+1) \\ 0 & \text{em contrário} \end{cases}$$

A (9) mostra que a moda é $y = 1$ e a expressão supra que o valor máximo de $f(y)$ é $\frac{k}{n}$.

Em conformidade com as (5), (6) e (7), tem-se :

$$\alpha_1 = \frac{n+1}{k+1}$$

$$\mu_2 = \frac{k(n+1)(n-k)}{(k+1)^2(k+2)}$$

$$\mu_3 = \frac{(n+1)(k-1)(2n-k+1)(n-k)}{(k+1)^3(k+2)(k+3)} < 0$$

Dos resultados acima segue que o gráfico representativo da função é da forma de J invertido.

f) - Se $n = k$ ($k > 1$) tem-se para a distribuição do extremo superior da graduatória amostral :

$$f(y) = \begin{cases} \frac{\binom{y-1}{k-1}}{\binom{n}{k}} & \text{se } y = k, k+1, \dots, n \\ 0 & \text{em qualquer outro caso} \end{cases}$$

A (9) mostra que a moda é $y = n$ e a expressão supra que o valor máximo de $f(y)$ é $\frac{k}{n}$. Em conformidade com as (5), (6) e (7) tem-se :

$$\alpha_1 = \frac{k(n+1)}{k-1}$$

$$\mu_2 = \frac{k(n+1)(n-k)}{(k+1)^2(k+2)}$$

$$\mu_3 = \frac{-k(n+1)(k-1)(2n-k+1)(n-k)}{(k+1)^3(k+2)(k+3)} < 0$$

Dos resultados acima segue que o gráfico representativo da função é da forma de J .

É fácil verificar que a atual distribuição coincide com a estudada em e) quando se submete esta a uma translação de amplitude $n-1$ seguida de uma rotação de π em torno de um eixo vertical levantado perpendicularmente ao eixo das abcissas no seu ponto n .

g) ~ A última propriedade salientada em f) relacionando as distribuições dos valores que ocupam o primeiro posto na co e contragraduatórias pode ser facilmente generalizada. Efetivamente, se na (1) substituirmos o parâmetro h por $k-h+1$, a distribuição assume a forma:

$$f_1(y) = \begin{cases} \frac{\binom{y-1}{k-k} \binom{n-y}{k-1}}{\binom{n}{k}} & \text{se } y = (k-k+1), \dots, (n-k+1) \\ 0 & \text{em contrário} \end{cases}$$

Comparando-se $f_1(y)$ com $f(y)$, é imediato verificar que :

$$f(k+\delta) = f_1(n-k-\delta+1) \quad (\delta = 0, 1, \dots, n-k)$$

De fato:

$$f(k+\delta) = \frac{\binom{k+\delta-1}{k-1} \binom{n-k-\delta}{k-k}}{\binom{n}{k}} = f_1(n-k-\delta+1)$$

Considerando-se, agora, que um indivíduo que ocupa o posto $k-h+1$ na cograduatória ocupa na contragraduatória o posto h , estabeleçamos, pois, genericamente, que o gráfico representativo da distribuição dos valores que ocupam o posto h na contragraduatória coincide com o representativo da distribuição daqueles que ocupam o mesmo posto h mas na cograduatória, quando se submete este último gráfico a uma translação de amplitude $n-2h+1$ seguida de uma rotação de π em torno de um eixo vertical levantado perpendicularmente ao eixo das abcissas no seu ponto $n-h+1$.

h)- Vamos supor aqui que h é o posto amostral mediano. Para evitar complicações formais, faremos $k=2p+1$ e, conseqüentemente $h = p+1$. A distribuição da mediana amostral é :

$$f(y) = \begin{cases} \frac{\binom{y-1}{p} \binom{n-y}{p}}{\binom{n}{2p+1}} & \text{se } y = p+1, \dots, n-p \\ 0 & \text{em contrário.} \end{cases}$$

Da (9) segue que a moda y deverá satisfazer as desigualdades :

$$\frac{n}{2} \leq y \leq \frac{n}{2} + 1$$

o que mostra que se n for par, haverá dois máximos e se n for ímpar, um único, dado por $\frac{n+1}{2}$.

Em conformidade com a (5), tem-se :

$$\alpha_1 = E(y) = \frac{n+1}{2}$$

oque, em palavras, pode-se traduzir, dizendo que a média dos valores dos postos populacionais que vêm a ocupar na amostra o posto mediano é igual ao posto mediano populacional.

Das (6), (7) e (8) deduzem-se :

$$\mu_2 = \frac{(n+1)(n-2p-1)}{4(2p+3)}$$

$$\mu_3 = 0$$

$$\mu_4 = \frac{(n+1)(n-2p-1)}{16(2p+3)(p+3)(2p+5)} (3n^2p - 6np^2 - 10p^2 - 12np - 27p + 6n^2 - 14)$$

$$\gamma_1 = 0$$

$$\gamma_2 = \frac{2p+3}{(n+1)(p+2)(2p+5)(n-2p-1)} (3n^2p - 6np^2 - 10p^2 - 12np - 27p + 6n^2 - 14)$$

Desta última expressão resulta que para $p \rightarrow \infty$ e conseqüentemente $n \rightarrow \infty$, a distribuição é me ocúr-tica ; com efeito, como nas variáveis n e p , a ordem do infinito do denominador de γ_2 é 4, tem-se:

$$\begin{aligned} \lim_{n,p \rightarrow \infty} r_2 &= \lim_{n,p \rightarrow \infty} \frac{(2p+3)(3n^2p - 6np^2)}{(n+1)(p+2)(2p+5)(n-2p-1)} = \\ &= \lim_{n,p \rightarrow \infty} \frac{(2 + \frac{3}{p})(3 - 6\frac{p}{n})}{(1 + \frac{1}{n})(1 + \frac{2}{p})(2 + \frac{5}{p})(1 - 2\frac{p}{n} - \frac{1}{n})} = 3 \end{aligned}$$

Resumindo, estabeleçamos: a distribuição dos valores que vêm a ocupar na amostra o posto mediano é simétrica em torno do ponto $\frac{n+1}{2}$, tem variância dada por $\frac{(n+1)(n-2p-1)}{4(2p+3)}$ e é mesocúrtica para amostras suficientemente grandes.

Para finalizar, demonstremos que para $p \rightarrow \infty$, a variável reduzida :

$$t = \frac{y - \frac{n+1}{2}}{\sqrt{\frac{(n+1)(n-2p-1)}{4(2p+3)}}}$$

é normal, com parâmetros 0 e 1.

Com efeito, pondo:

$$x = t \sqrt{\frac{(n+1)(2n-p-1)}{4(2p+3)}}$$

tem-se:

$$P\left(x + \frac{n+1}{2}\right) = P_y = \frac{\binom{\frac{n+2x-1}{2}}{p} \binom{\frac{n-2x-1}{2}}{p}}{\binom{n}{2p+1}}$$

Ora, o valor da probabilidade máxima sendo dada por:

$$P\left(\frac{n+1}{2}\right) = P_0 = \frac{\binom{\frac{n-1}{2}}{p} \binom{\frac{n-1}{2}}{p}}{\binom{n}{2p+1}}$$

pode-se escrever :

$$\frac{P_{2x}}{P_0} = \frac{\binom{\frac{n+2x-1}{2}}{p} \binom{\frac{n-2x-1}{2}}{p}}{\binom{\frac{n-1}{2}}{p} \binom{\frac{n-1}{2}}{p}} =$$

$$= \frac{\left(\frac{n+2x-1}{2}\right)! \left(\frac{n-2x-1}{2}\right)! \left(\frac{n-2p-1}{2}\right)! \left(\frac{n-2p-1}{2}\right)!}{\left(\frac{n+2x-1-2p}{2}\right)! \left(\frac{n-2x-1-2p}{2}\right)! \left(\frac{n-1}{2}\right)! \left(\frac{n-1}{2}\right)!}$$

do que segue, por utilização da fórmula de Stirling

$$n! \approx n^n \cdot e^{-n} \sqrt{2\pi n}$$

$$\frac{P_{2x}}{P_0} \approx \frac{\left(\frac{n+2x-1}{2}\right)^{\frac{n+2x}{2}} \left(\frac{n-2x-1}{2}\right)^{\frac{n-2x}{2}} \left(\frac{n-2p-1}{2}\right)^{n-2p}}{\left(\frac{n+2x-1-2p}{2}\right)^{\frac{n+2x-2p}{2}} \left(\frac{n-2x-1-2p}{2}\right)^{\frac{n-2x-2p}{2}} \left(\frac{n-1}{2}\right)^n}$$

$$= \frac{\left(1 + \frac{2x}{n-1}\right)^{\frac{n}{2}+x} \left(1 - \frac{2x}{n-1}\right)^{\frac{n}{2}-x} \left(\frac{n-1}{2}\right)^{\frac{n}{2}+x} \left(\frac{n-1}{2}\right)^{\frac{n}{2}-x} \left(\frac{n-2p-1}{2}\right)^{n-2p}}{\left(1 + \frac{2x}{n-2p-1}\right)^{\frac{n}{2}+x-p} \left(\frac{n-2p-1}{2}\right)^{\frac{n}{2}+x-p} \left(1 - \frac{2x}{n-2p-1}\right)^{\frac{n}{2}-x-p} \left(\frac{n-2p-1}{2}\right)^{\frac{n}{2}-x-p} \left(\frac{n-1}{2}\right)^n}$$

$$= \frac{\left(1 + \frac{2x}{n-1}\right)^{\frac{n}{2}+x} \left(1 - \frac{2x}{n-1}\right)^{\frac{n}{2}-x}}{\left(1 + \frac{2x}{n-2p-1}\right)^{\frac{n}{2}+x-p} \left(1 - \frac{2x}{n-2p-1}\right)^{\frac{n}{2}-x-p}}$$

Assim sendo, tem-se :

$$\begin{aligned} \log \frac{P_1}{P_0} &\approx \left(\frac{n}{2} + x\right) \log \left(1 + \frac{2x}{n-1}\right) + \left(\frac{n}{2} - x\right) \log \left(1 - \frac{2x}{n-1}\right) - \\ &- \left(\frac{n}{2} + x - p\right) \log \left(1 + \frac{2x}{n-2p-1}\right) - \left(\frac{n}{2} - x - p\right) \log \left(1 - \frac{2x}{n-2p-1}\right) = \\ &= \left(\frac{n}{2} + x\right) \left[\frac{2x}{n-1} - \frac{4x^2}{2(n-1)^2} + \dots \right] - \left(\frac{n}{2} - x\right) \left[\frac{2x}{n-1} + \frac{4x^2}{2(n-1)^2} + \dots \right] - \\ &- \left(\frac{n}{2} + x - p\right) \left[\frac{2x}{n-2p-1} - \frac{4x^2}{2(n-2p-1)^2} + \dots \right] + \left(\frac{n}{2} - x - p\right) \left[\frac{2x}{n-2p-1} + \frac{4x^2}{2(n-2p-1)^2} + \dots \right] = \\ &= -n \left[\frac{4x^2}{2(n-1)^2} + \frac{16x^4}{4(n-1)^4} + \dots \right] + n \left[\frac{4x^2}{2(n-2p-1)^2} + \frac{16x^4}{4(n-2p-1)^4} + \dots \right] + \\ &+ 2x \left[\frac{2x}{n-1} + \frac{8x^3}{3(n-1)^3} + \dots \right] - 2x \left[\frac{2x}{n-2p-1} + \frac{8x^3}{3(n-2p-1)^3} + \dots \right] - \\ &- 2p \left[\frac{4x^2}{2(n-2p-1)^2} + \frac{16x^4}{4(n-2p-1)^4} + \dots \right] \end{aligned}$$

Considerando que para n e p tendentes ao infinito, x é um infinito de ordem $1/2$, então :

$$\begin{aligned} \lim_{n,p \rightarrow \infty} \log \frac{P_1}{P_0} &= \lim_{n,p \rightarrow \infty} \frac{4x^2(n+1)(n-2p-1)}{4(2p+3)} \left\{ -\frac{n}{2(n-1)^2} + \frac{n}{2(n-2p-1)^2} + \frac{1}{n-1} \right. \\ &\quad \left. - \frac{1}{n-2p-1} - \frac{p}{(n-2p-1)^2} \right\} \\ &= x^2 \left(\frac{1}{4} - \frac{1}{4} - \frac{1}{2} - \frac{1}{2} + \frac{1}{2} \right) = -\frac{x^2}{2} \end{aligned}$$

donde, finalmente:

$$\lim_{n,p \rightarrow \infty} \frac{P_1}{P_0} = e^{-x^2/2} \quad \text{c.s.q.d.}$$

4 - A grande importância prática da nossa distribuição consiste em fornecer a base para inferências lí-

citadas quanto às possíveis posições na população de um elemento que na amostra ocupa o posto h . Assim, por exemplo, podemos afirmar que a observação de posto h na amostra tem probabilidade $1/n$ de ocupar um dos postos populacionais do intervalo

$$h \text{ --- } n-k+h$$

pois este é precisamente o campo de definição de y .

Este intervalo, em geral, não é desejável, devido a sua amplitude, com o que somos levados à procura de um novo, mais restrito, constituído por um conjunto de valores y tendo um total de probabilidade de ocupar na amostra o posto h dado por P .

Por seu turno, a solução deste último problema está na dependência da função cumulativa de y , isto é, do valor de :

$$\frac{1}{\binom{n}{h}} \sum_{y=h}^{h+m} \binom{y-1}{h-1} \binom{n-y}{n-h} \quad (m=0,1,\dots,n-k) \quad (11)$$

Como sempre acontece com as distribuições discretas, também aqui não é possível obter uma expressão de maneio cómodo para a (11), razão pela qual, escolhido o tipo de intervalo - o que naturalmente depende de h^* - teremos de

* Devido a forma da distribuição, parece-nos lógico que segundo h esteja mais próximo de $1, k+1/2, k$, as regiões de rejeição sejam, respectivamente, a cauda superior, as caudas superior e inferior e a cauda inferior.

calcular, uma por uma, as probabilidades dos y da região a ser subtraída de $h = n - k + h$, até que sua soma faça um total o mais aproximado possível de $1 - P$.

Com a finalidade única de concretizar estas considerações, vamos abordar aqui a hipótese de h ser o posto amostral mediano e supor novamente $k = 2p + 1$. Tendo em vista que no caso em apreço, a distribuição é simétrica em torno do ponto médio do intervalo $p + 1 \rightarrow n - p$ e isto é, de $n + 1/2$, parece-nos lógico que a região de rejeição seja dada pelas duas caudas da distribuição. Adotado este ponto de vista, o nosso problema resume-se na determinação de um valor de m tal que :

$$\sum_{j=p+1}^{(p+1)+m} \frac{\binom{y-1}{p} \binom{m-y}{p}}{\binom{n}{2p+1}} = \sum_{j=0}^m \frac{\binom{j+p}{p} \binom{n-j-p-1}{p}}{\binom{n}{2p+1}} \approx \frac{1-P}{2}$$

Para o particular caso $n = 40$, $p = 10$ e $P = 0,95$, teríamos de encontrar um valor de m tal que:

$$E_m = \sum_{j=0}^m \frac{\binom{j+10}{10} \binom{29-j}{10}}{\binom{40}{21}} \approx 0,025$$

Ora, tem-se:

$$E_0 = 0,015\% \quad E_1 = 0,015\% + 0,11\% = 0,125\%$$

$$E_2 = 0,125\% + 0,42\% = 0,545\% \quad E_3 = 0,545\% + 1,16\% = 1,705\%$$

$$E_4 = 1,705\% + 2,5\% = 4,205\%$$

Dêstes resultados seguen $m = 3$ e, portanto, o intervalo de confiança 95% dado por 15 — 26

Como vemos, a construção do intervalo procurado exige, mesmo neste exemplo simples, não pequeno trabalho de cálculo. Em casos ainda mais complicados, para obviar este inconveniente, poderíamos, naturalmente, lembrar de expedientes outros, como sejam tabelas de fatoriais, emprego de logarítmos, fórmulas de Stirling, etc...

Voltando ao caso do posto mediano, além da utilização da distribuição limitante, poderíamos ainda lançar mão da propriedade empírica, segundo a qual em distribuições simétricas, a amplitude $\mu \pm 2\sigma$ abrange cerca de 95% das observações. Este último expediente nos levaria ao intervalo :

$$\frac{n+1}{2} \pm \sqrt{\frac{(n+1)(n-2p+1)}{2p-1}} \quad \text{ou} \quad \frac{n+1}{2} \pm \sqrt{\frac{(n+1)(n-2p+1)}{2p-1}}$$

que, no exemplo numérico apresentado torna-se $20,5 \pm \sqrt{\frac{21 \cdot 19}{23}}$ ou seja 14,7 — 26,3 que, neste caso, coincide com o intervalo exato.

5 - Na epígrafe 1 salientamos que devido à correspondência biunívoca entre os números naturais e a grandeza postos, a solução do problema objeto dêste trabalho podia ser interpretada como dando a distribuição da variável aleatória " valor que ocupa o posto h " quando se amostra, de forma casual, sem reposição, uma população cons-

tituída dos n primeiros números naturais. Esta hipótese da população ser dada pelos valores i ($i = 1, 2, \dots, n$) pode ser substituída pela mais geral de que a população consiste dos elementos $a+1r, a+2r, \dots, a+nr$ de uma progressão aritmética.

De fato, tem-se imediatamente para distribuição do valor $a+ry$ que ocupa o posto amostral h :

$$P(a+ry) = P(y) = \frac{\binom{y-1}{h-1} \binom{n-y}{k-h}}{\binom{n}{k}} \quad y=h, h+1, \dots, n-k$$

Dêste resultado segue que os momentos centrais da referida variável "valor que ocupa o posto h " são dados pelas (6), (7) e (8) multiplicadas respectivamente por r^2, r^3 e r^4 , ao passo que se tem para a média dela :

$$E(a+ry) = a + \frac{h(n+1)}{k+1} \cdot r$$

e para sua moda o valor (ou valores) que se obtém ao substituir em $a+ry$, y pelo inteiro (ou inteiros) definido pelas desigualdades (9).

Não é, entretanto, a generalização acima, a que nos parece mais importante, mas sim a que resulta da resolução do problema abaixo, generalização direta do nosso primitivo problema.

" Uma população contém n indivíduos numerados de 1 a n , segundo um certo critério. Toma-se uma amostra casual de k ($k \leq n$) indivíduos, sem reposição,

e cogradua-se segundo o posto por eles ocupado na população. Pergunta-se : qual a probabilidade de que os indivíduos que na população ocupam os postos $y_1 < y_2 < \dots < y_r$ ($r \leq k$), virem a ocupar na amostra os postos respectivos $1, h_1, \dots, h_r$.

Para que o evento se verifique é necessário e suficiente que na amostra apareçam $h_1 - 1$ indivíduos com postos na população menor do que y_1 , os indivíduos de posto y_1, y_2, \dots, y_r , $(h_j - h_{j-1} - 1)$ com $j=2, 3, \dots, r$ indivíduos com postos, na população, compreendidos entre y_{j-1} e y_j e finalmente, $k - (h_1 - 1) - \sum_{j=2}^r (h_j - h_{j-1} - 1) - r = k - h_r$ indivíduos de postos, na população, maiores do que y_r .

Ora, $h_1 - 1$ indivíduos com postos, na população, menores do que y_1 podendo ser escolhidos de $\binom{y_1 - 1}{h_1 - 1}$ diferentes maneiras, $h_j - h_{j-1} - 1$ com postos, na população, compreendidos entre y_{j-1} e y_j podendo ser escolhidos de $\binom{y_j - y_{j-1} - 1}{h_j - h_{j-1} - 1}$ $j=2, 3, \dots, r$ diferentes maneiras, e os $k - h_r$ de postos, na população, maiores do que y_r de $\binom{n - y_r}{k - h_r}$ diferentes maneiras, a probabilidade procurada vale, desde que existem $\binom{n}{k}$ amostras possíveis :

$$f(y_1, \dots, y_r) = \frac{\binom{y_1 - 1}{h_1 - 1}}{\binom{n}{k}} \left\{ \prod_{j=2}^r \binom{y_j - y_{j-1} - 1}{h_j - h_{j-1} - 1} \right\} \binom{n - y_r}{k - h_r} \quad (12)$$

com $y_\beta = h_\beta, h_\beta + 1, \dots, y_{\beta+1} - (h_{\beta+1} - h_\beta)$ $\beta = 1, 2, \dots, r-1$

$y_r = h_r, h_r + 1, \dots, n - k - h_r$

É fácil demonstrar que :

$$I = \sum_{\substack{n-k+h_n \\ y_n = h_n}} \sum_{\substack{y_n = (h_n - h_{n-1}) \\ y_{n-1} = h_{n-1}}} \sum_{y_1 = h_1}^{y_n - (h_2 - h_1)} f(y_2, \dots, y_n) = 1$$

De fato :

$$I = \frac{1}{\binom{n}{k}} \sum_{y_n = h_n}^{n-k+h_n} \binom{n-y_n}{k-h_n} \sum_{y_{n-1} = h_{n-1}}^{y_n - (h_n - h_{n-1})} \binom{y_n - y_{n-1} - 1}{h_n - h_{n-1} - 1} \dots$$

$$\dots \sum_{y_2 = h_2}^{y_3 - (h_3 - h_2)} \binom{y_3 - y_2 - 1}{h_3 - h_2 - 1} \sum_{y_1 = h_1}^{y_2 - (h_2 - h_1)} \binom{y_2 - y_1 - 1}{h_2 - h_1 - 1} \binom{y_2 - 1}{h_2 - 1}$$

Aplicando sucessivamente a (1), vem :

$$I = \frac{1}{\binom{n}{k}} \sum_{y_n = h_n}^{n-k+h_n} \binom{n-y_n}{k-h_n} \sum_{y_{n-1} = h_{n-1}}^{y_n - (h_n - h_{n-1})} \binom{y_n - y_{n-1} - 1}{h_n - h_{n-1} - 1} \dots$$

$$\dots \sum_{y_2 = h_2}^{y_3 - (h_3 - h_2)} \binom{y_3 - y_2 - 1}{h_3 - h_2 - 1} \binom{y_2 - 1}{h_2 - 1} =$$

=

$$= \frac{1}{\binom{n}{k}} \sum_{y_n = h_n}^{n-k+h_n} \binom{n-y_n}{k-h_n} \binom{y_n - 1}{h_n - 1} = 1 \quad \text{c. s. q. d.}$$

Da (12) segue :

$$f_{2\dots p}(y_2, y_3, \dots, y_p) = \sum_{y_p = k_{p-1}}^{y_2 - (k_p - k_{p-1})} \dots \sum_{y_2 = k_1}^{y_2 - (k_2 - k_1)} f(y_1, y_2, \dots, y_n) =$$

$$= \frac{\binom{n - y_n}{k - k_n} \left\{ \prod_{j=p+1}^n \binom{y_j - y_{j-1} - 1}{k_j - k_{j-1} - 1} \right\} \binom{y_n - 1}{k_p - 1}}{\binom{n}{k}} \quad (13)$$

o que mostra que a marginal das $r = p + 1$ últimas variáveis y_r, y_{r-1}, \dots, y_p é de forma análoga a da (12).

Da (13) segue, por seu turno, para a condicional das $p - 1$ primeiras variáveis :

$$f(y_1, y_2, \dots, y_{p-1} / y_p, \dots, y_n) = \frac{f(y_1, \dots, y_{p-1}, y_p, \dots, y_n)}{f_{2\dots p}(y_2, \dots, y_p)} =$$

$$= \frac{\binom{y_1 - 1}{k_1 - 1} \prod_{j=2}^p \binom{y_j - y_{j-1} - 1}{k_j - k_{j-1} - 1}}{\binom{y_p - 1}{k_p - 1}} \quad (14)$$

Como caso particular da (12), temos, pondo $r = 2$, a distribuição bidimensional :

$$f(y_2, y_1) = \frac{\binom{y_1 - 1}{k_1 - 1} \binom{y_2 - y_1 - 1}{k_2 - k_1 - 1} \binom{n - y_2}{k - k_2}}{\binom{n}{k}} \quad (15)$$

$$y_1 = k_1, k_1 + 1, \dots, y_2 - k_2 + k_1$$

$$y_2 = k_2, k_2 + 1, \dots, n - k_2 + k_1$$

do que segue :

$$f_2(y_2) = \frac{\binom{n - y_2}{k - k_2} \binom{y_2 - 1}{k_2 - 1}}{\binom{n}{k}}$$

com o que recaímos, como era de se esperar, na distribuição originária.

Tem-se, também :

$$f_2(y_1/y_2) = \frac{\binom{y_1-1}{h_1-1} \binom{y_2-y_1-1}{h_2-h_1-1}}{\binom{y_2-1}{h_2-1}}$$

Fazendo na (15) $h_1 = 1$ e $h_2 = k$, vem :

$$f(y_1, y_2) = \frac{\binom{y_2-y_1-1}{k-2}}{\binom{n}{k}} \quad \begin{array}{l} y_1 = 1, 2, \dots, y_2-k+1 \\ y_2 = k, k+1, \dots, n \end{array} \quad (16)$$

que é a distribuição conjunta dos valores extremos amostrais.

Da (16) segue, por seu turno, para as marginais :

$$f_2(y_2) = \frac{\binom{y_2-1}{k-1}}{\binom{n}{k}} \quad y_2 = k, k+1, \dots, n$$

$$f_1(y_1) = \frac{\binom{n-y_1}{k-1}}{\binom{n}{k}} \quad y_1 = 1, 2, \dots, n-k+1$$

com o que recaímos nas distribuições consideradas em f) e e) da epígrafe 4.

Ainda da (16) segue para a distribuição da amplitude t_1 de variação amostral :

$$f(t_1) = \frac{\binom{n-t_1}{k-2} \binom{t_1-1}{k-2}}{\binom{n}{k}} \quad t_1 = k-1, k, \dots, n-1.$$

que é a própria (1) quando se faz $h = k - 1$.

6 - Do ponto de vista teórico, dois importantes papéis podem ser atribuídos à nossa distribuição, visto como demonstraremos os dois fatos abaixo:

1^o) - através dela é possível determinar a distribuição da variável aleatória "área sob a função de frequência de uma população finita à esquerda do valor que numa amostra, sem reposição, ocupa o posto h " e provar que tal distribuição independe da função de frequência da população;

2^o) - demonstrar, no caso bastante geral de uma população enumerável, o teorema básico das provas não paramétricas, ou seja, determinar a distribuição da variável aleatória "área sob a função de frequência de uma população à esquerda de um valor que numa amostra simples ocupa o posto h ".

Isto é conseguido da seguinte maneira: Vimos, pela resolução do nosso problema que um valor ocupando na população postos variáveis de h até $n - k + h$ podem ocupar na amostra o posto h com probabilidades dadas pela (1).

Nossa única suposição ao estabelecer tais probabilidades foi a de admitir uma graduação na população finita amostrada, nenhuma hipótese tendo sido formulada quanto à natureza desta população, isto é, quanto a sua densidade de probabilidade que simbolisaremos por $\psi(x)$. Considerando, agora, que a área sob a $\psi(x)$ à esquerda de um valor que na amostra ocupa o posto h depende e só depende do posto por ele ocupado na população, então, a distribuição de A independe, como afirmamos, de $\psi(x)$. Além disto, subsiste a seguinte correspondência:

Valores de A	Probabilidades
$\frac{k-1}{n}$	$f(k)$
$\frac{k}{n}$	$f(k+1)$
\vdots	
$\frac{n-k+k-1}{n}$	$f(n-k+k)$

resultado que mostra que a distribuição de A pode ser obtida da (1) quando se submete esta à transformação

$$y = nA + 1 \quad (17)$$

Tem-se, então :

$$P(A) = \frac{\binom{nA}{k-1} \binom{n-nA-1}{k-k}}{\binom{n}{k}}$$

$$A = \frac{k-1}{n}, \dots, \frac{n-k+k-1}{n}$$

Suponhamos, a seguir, $n \rightarrow \infty$, com o que a amostragem sem reposição transforma numa amostragem simples. No caso em apreço, A varia entre 0 e 1 e, tem-se, notando que o jacobiano da (17) vale n :

$$\lim_{n \rightarrow \infty} p(A) = y(A) = \lim_{n \rightarrow \infty} \frac{\binom{nA}{k-1} \binom{n-nA-1}{k-k}}{\binom{n}{k}} \cdot n$$

ou seja :

$$g(A) = \frac{k! A^{k-1} (1-A)^{k-k}}{(k-1)! (k-k)!}$$

que é a distribuição $\beta(k, k-k+1)$, c. s. q. d.

7 - Recordando que a variável y se identifica a $nA + 1$, então, simbolizando por A_1 ($1=1,2,\dots,r$) a área sob a $\psi(x)$ à esquerda do valor que na amostra ocupa o posto h_1 , estabeleçamos, como generalização:

$$p(A_1, A_2, \dots, A_r) = \frac{1}{\binom{n}{k}} \binom{nA_1}{h_1-1} \left[\prod_{j=2}^r \binom{nA_j - nA_{j-1} - 1}{h_j - h_{j-1} - 1} \right] \binom{n - nA_r - 1}{k - h_r}$$

$$A_\beta = \frac{h_\beta - 1}{n}, \frac{h_\beta}{n}, \dots, \frac{A_{\beta+1} - (h_\beta - h_\beta) - 1}{n} \quad \beta = 1, 2, \dots, r-1$$

$$A_r = \frac{h_r - 1}{n}, \dots, \frac{n - k + h_r - 1}{n}$$

Notando que o jacobiano é n^r , tem-se, para $n \rightarrow \infty$:

$$\lim_{n \rightarrow \infty} p(A_1, A_2, \dots, A_r) = g(A_1, A_2, \dots, A_r) =$$

$$= \frac{n^r \binom{nA_1}{h_1-1} \left[\prod_{j=2}^r \binom{nA_j - nA_{j-1} - 1}{h_j - h_{j-1} - 1} \right] \binom{n - nA_r - 1}{k - h_r}}{\binom{n}{k}}$$

$$= \frac{k! A_1^{k-1} \prod_{j=2}^r (A_j - A_{j-1})^{h_j - h_{j-1} - 1} (1 - A_r)^{k - h_r}}{(h_1 - 1)! \prod_{j=2}^r (h_j - h_{j-1} - 1)! (k - h_r)!} \quad (18)$$

$$0 \leq A_1 < A_2 < \dots < A_r \leq 1$$

O segundo membro da (18) é a distribuição conjunta das r áreas sob a $\psi(x)$ à esquerda dos valores que ocupam na amostra os postos h_1 numa amostragem casual simples.

Em particular, se $r = k$ e conseqüentemente $h_1 = 1$, tem-se para a distribuição das k áreas sob a $\psi(x)$ à esquerda dos k valores de uma amostra casual simples :

$$g(A_1, A_2, \dots, A_r) = k!$$

Ainda como caso particular da (18), tem-se, supondo $r = 2$:

$$g(A_1, A_2) = \frac{k! A_1^{h_1-1} (A_2 - A_1)^{h_2 - h_1 - 1} (1 - A_2)^{k - h_2}}{(h_1 - 1)! (h_2 - h_1 - 1)! (k - h_2)!}$$

$$0 < A_1 < 1 - A_2 < 1$$

da qual resulta, para densidade de probabilidade da variável $t = A_2 - A_1$, isto é, a área sob a $\psi(x)$ compreendida entre os valores que na amostra ocupam os postos h_2 e h_1 :

$$u(t) = \frac{k! t^{h_2 - h_1 - 1} (1 - t)^{k - h_2 - h_1}}{(h_2 - h_1 - 1)! (k - h_2 + h_1)!} \quad 0 < t < 1$$

que é novamente uma função do tipo beta, com parâmetros $h_2 - h_1$ e $k - h_2 + h_1 + 1$

Para finalizar, diremos que a importância dos resultados conseguidos nesta epígrafe na teoria dos testes não paramétricos é de tal ordem, que é desnecessário encarecê-la.

SUMÁRIO E CONCLUSÃO

I - Procurando determinar a probabilidade de um indivíduo que numa graduação populacional com n constituintes ocupa o posto y vir a ocupar na amostra o posto h , em uma amostragem, sem reposição, de tamanho k , fomos conduzidos à expressão :

$$f(y) = \frac{\binom{y-1}{h-1} \binom{n-y}{k-h}}{\binom{n}{k}} \quad (1)$$

II - Pela messe de resultados que fornece, a (1) deve ser enquadrada entre as mais importantes distribuições discretas de probabilidades.

III - A (1) fornece uma demonstração estatística elementar da fórmula do cálculo combinatório :

$$\sum_{z=0}^{n-k} \binom{z+h-1}{h-1} \binom{n-z}{k-h} = \binom{n}{k}$$

IV - A maneira mais adequada de conseguir as expressões dos momentos da (1) é fazê-lo por via dos momentos fatoriais da variável $n - y - k + h$, obtendo-se, então as fórmulas (3) e (4) do texto.

V - A distribuição (1) admite dois ou um único máximo, segundo $n(h-1)$ for ou não múltiplo de $k-1$.

VI - A (1) contém, como casos particulares, as dis-

tribuição ξ e a distribuição uniforme.

VII - Foram deduzidas expressões gerais para a soma das potências e para os momentos centrados dos n primeiros números naturais.

VIII - Como outros interessantes casos particulares da (1) foram estabelecidas as distribuições dos postos populacionais que ocupam na graduatória amostral o extremo inferior, o extremo superior e o posto mediano.

IX - O gráfico representativo dos valores que ocupam o posto h na contragraduatória amostral coincide com o representativo da distribuição dos valores que ocupam o mesmo posto h mas na cograduatória amostral, desde que se submeta este a uma translação de amplitude $n-2h+1$ seguida de uma rotação de π em torno de um eixo vertical levantado no ponto $n-h+1$ do eixo das abcissas.

X - A distribuição limitante da distribuição dos postos populacionais que na amostra ocupam o posto mediano é normal.

XI - A grande importância prática da nossa distribuição consiste em fornecer a base para a construção de intervalos de confiança dando, com probabilidades prefixadas, a posição na população de um indivíduo que na amostra ocupa o posto h .

XII - Como generalização da (1) foi estabelecida, para o caso r -dimensional, a fórmula (12), a distribuição marginal das $r-p+1$ últimas variáveis desta distribui-

ção constando da (13) e a condicional das $p-1$ primeiras variáveis aparecendo na (14) .

XIII - Como caso particular da (12) obteve-se a distribuição bidimensional (16) por via da qual deduziu-se a distribuição da amplitude de variação amostral.

XIV - Partindo da (1) foi possível demonstrar, no caso bastante geral de uma população enumerável, a proposição básica da teoria dos testes não paramétricos.

XV - Outros importantes resultados para esta última teoria foram estabelecidos na epígrafe 7 .

BIBLIOGRAFIA

- Berzolari, L. , Vivanti, G. e Gigli, O. -
Enciclopedia delle matematiche elementari, Vol. I-Parte II. Milano, Giucoco Hoepli, 1931.
- Cramér, H. -
Mathematical Methods of Statistics. Princeton, Princeton University Press, 1946.
- Duarte, G. G. -
Contribuição para o estudo dos momentos fatoriais. Tese. São Paulo, Grafico Omega, 1950.
- Kendall, M. G. -
The advanced theory of Statistics. Vol I . London. J. Lippincott Company, 1942.
- Wood, A. M. -
Introduction to the theory of Statistics. New York. Mc-Graw-Hill Book Company, Inc. , 1950.
- Wilks, S. S. -
Mathematical Statistics. Princeton. Princeton University Press, 1946.