

AS DISTRIBUIÇÕES DO ACASO

F. G. Brieger

Seção de Genética
Escola Superior de Agricultura
"Lutz de Queiroz", Universidade
de S. Paulo

INDICE

Prefácio	322
1 — <i>Introdução</i>	323
2 — <i>A análise do erro standard</i>	328
3 — <i>A definição matemática dos quatro tipos de distribuição de acaso</i> ..	331
a) <i>As fórmulas das distribuições</i>	331
b) <i>Cálculo das ordenadas e áreas</i>	334
c) <i>Forma das curvas</i>	334
d) <i>Comparação matemática das quatro distribuições</i>	335
4 — <i>A transformação matemática das quatro distribuições</i>	341
a) <i>Noções matemáticas</i>	341
b) <i>Transformação geral da parte exponencial</i>	343
c) <i>Transformação da equação de Fisher para a de Pearson</i>	345
d) <i>Transformação da equação de Pearson para a de Gauss</i>	348
e) <i>Transformação da equação de Fisher para a de Student</i> ...	349
f) <i>Transformação da equação de Student para a de Gauss</i>	350
5 — <i>A forma das distribuições em dependência do grau de liberdade</i> ...	351
a) <i>As distribuições de Student</i>	351
b) <i>As distribuições de Pearson</i>	351
c) <i>Distribuições de Fisher</i>	353
6 — <i>Relação entre as distribuições e os testes derivados</i>	358
7 — <i>Considerações finais</i>	361
<i>Resumo</i>	372
<i>Abstrat</i>	377
<i>Referências</i>	382

PREFÁCIO

O presente trabalho trata das bases fundamentais da teoria da variação de acaso, e assim tornou-se inevitável entrar em discussões que requerem algumas noções de matemática. Porém para que o trabalho seja acessível para um maior número de leitores, limitei a aplicação de termos matemáticos ao mínimo necessário.

Um conhecimento das definições matemáticas e das relações entre as diferentes distribuições do acaso é indispensável para os técnicos que querem especializar-se profundamente na matéria e idealizar novos métodos da análise. Mas devemos acentuar que estes conhecimentos profundos são dispensáveis para a grande maioria daqueles que pretendem apenas aplicar os métodos da análise estatística em seus estudos experimentais. Para isto basta de modo geral seguir as normas estabelecidas e os testes recomendados por técnicos tão reconhecidos como KARL PEARSON, STUDENT, R. A. FISHER e outros.

Uma situação análoga encontramos em quase todos os trabalhos experimentais. Continuamente o especialista em determinado assunto precisa aplicar métodos de outros ramos, também especializados. Assim, o citologista, colorindo as suas lâminas com violeta genciana, com hematoxilina ou outros corantes, não precisa saber a fórmula ou a origem desses compostos, nem tampouco detalhes químicos do processo da coloração, detalhes estes muitas vezes desconhecidos. Um químico, que quer apenas determinar a porcentagem de açúcar em extratos ou valor do pH em soluções, aplica em geral métodos químicos ou físicos, dispensando os conhecimentos da teoria ótica da polarização da luz, da teoria fisico-química da mudança de cor dos indicadores e dos problemas fundamentais da teoria da condutibilidade elétrica e da ionização.

De um modo geral, um especialista deve saber todas as teorias da sua especialidade, mas quando precisa aplicar métodos de outras especialidades, ele pode e deve confiar na competência daqueles que elaboraram os respectivos métodos.

O presente trabalho é portanto destinado para os estatísticos e não para os técnicos que querem apenas aplicar os métodos de análise na experimentação geral.

1 — INTRODUÇÃO (1)

Nos últimos quarenta anos a aplicação de novos métodos de análise estatística recebeu um grande impulso e também uma extensão pelas descobertas da escola londrina de KARL PEARSON, STUDENT e R. A. FISHER. Notamos, nos trabalhos mais recentes de estatística, nitidamente dois grupos: um que continua a dar quase que exclusivamente os métodos anteriores às descobertas de PEARSON (1900) e outro que apresenta os métodos modernos sem porém explicar a história de sua descoberta e a sua relação com os métodos da estatística anterior, nem mesmo as suas relações entre si. É o fim desta publicação mostrar a relação entre alguns dos diferentes métodos, modernos e velhos, e dos testes deles derivados. Estes últimos, especialmente, são em geral explicados numa forma matemática tão alterada, que não é fácil entender a sua derivação nem as suas relações lógicas com os demais testes.

O fim de toda a análise estatística é mostrar se valores observados correspondem aos valores que podem ser calculados de acordo com uma teoria ou hipótese. Devemos de um modo geral decidir estatisticamente se a diferença entre estes

(1) Parece-me indicado no início do trabalho dar algumas explicações sobre os símbolos que serão usados, pois até hoje não existe nenhum acordo científico sobre o assunto.

Valores variáveis indicaremos com qualquer letra: v ou a , etc.. As médias serão caracterizadas por traços horizontais sobre as letras: \bar{v} \bar{a} , etc.. Quando forem determinadas as médias parciais e gerais, as primeiras serão indicadas por um traço: \bar{v} e as médias gerais por dois traços $\bar{\bar{v}}$.

A dificuldade maior encontramos quando queremos definir os símbolos relativos à frequência das variáveis. Temos em geral que distinguir o número total (N) de variáveis e o número de variáveis livres, que em outras publicações indiquei sempre pelas letras duplas nf . Mas para simplificar a impressão, especialmente quando temos ainda que usar sufixos com nf_1 , nf_a , etc., usarei em geral neste trabalho apenas a letra n para indicar o grau de liberdade, recorrendo apenas à forma nf quando há perigo de mal entendidos.

Nas referências às fórmulas da estatística clássica, isto é, especialmente na introdução, usarei a letra n , sem distinguir claramente entre o número total de variáveis e o seu grau de liberdade. Uma vez que estas discussões se referem apenas a amostras grandes, estes dois conceitos podem, sem prejuízo, ser confundidos.

dois valores, observado e esperado, medido pelo valor do respectivo erro standard, pode ou não ser atribuída ao acaso.

Denominando o valor observado por v , o valor ideal por v^0 , ou quando êle é representado por uma média aritmética por \bar{v} , e usando finalmente a letra grega σ para o erro standard da variável v , temos assim a fórmula do desvio relativo simples:

$$D = \frac{v - v_0}{\sigma} \text{ ou } = \frac{v - \bar{v}}{\sigma} \text{ ----- 1-1}$$

Supoz-se que a variação dêste termo seguiria sempre a fórmula matemática que caracteriza a distribuição normal ou de GAUSS, descoberta independentemente por DE MOIVRE, LA PLACE e GAUSS. É o mérito dos autores já mencionados ter demonstrado que a chamada distribuição "normal" é apenas um caso, muito especial, entre um número grande de outras distribuições do acaso. Para melhor demonstrar em que consiste esta particularidade da distribuição de GAUSS, vamos dar em primeiro lugar rapidamente a derivação da fórmula matemática da mesma.

Supomos que temos um aparelho de acaso de GALTON, bem construído, no qual um número de bolas de aço, caindo verticalmente, encontra no seu caminho n pregos arranjados de tal modo que a probabilidade de ricochetear para a direita é p e para a esquerda q . Podemos então determinar as frequências de tôdas as combinações de desvios para a direita ou para a esquerda, calculando os termos de um binômio:

$$(p + q)^n = q^n + \frac{n!}{(n-1)!} \cdot p \cdot q^{n-1} + \frac{n!}{2!(n-2)!} \cdot p^2 \cdot q^{n-2} \dots \text{---} \\ \text{----- 1-2}$$

De um modo geral podemos exprimir a frequência y de um termo m nesta série, pela equação:

$$y = \frac{n!}{m!(n-m)!} \cdot p^m \cdot q^{n-m} \text{----- 1-3}$$

sendo a frequência máxima dos termos centrais das séries binomiais:

$$y(m) = \frac{n!}{np! nq!} \cdot p^{np} \cdot q^{nq} \quad \text{----- 1-4}$$

Podemos, de outro lado, dar a frequência de qualquer termo m , usando para a sua caracterização a sua distância x do termo central, pelas fórmulas:

$$m = np - x$$

$$n - m = n - (np - x) = nq + x$$

sendo conseqüentemente

$$y = \frac{n!}{(np-x)(nq+x)} \cdot p^{np-x} \cdot q^{nq+x} \quad \text{----- 1-5}$$

Aplicando agora alguns princípios do cálculo infinitesimal, podemos transformar esta equação geral, para obter a solução do caso especial no qual n atinge o valor infinito. Os gráficos na fig. 1 explicam porem ainda melhor o que acontece quando o valor do expoente n do binômio cresce. Estão reproduzidos quatro histogramas, que correspondem a valores escolhidos de $n = 4, 16, 36$ e 64 . O erro standard, calculado pela fórmula.

$$\sigma = \pm \sqrt{p \cdot (1-p)n}$$

tem então os valores: $\pm 1, \pm 2, \pm 3$ e ± 4 . As frequências de cada classe são sempre indicadas acima de cada coluna do histograma. Nota-se que com o crescimento de valor de n , aumenta o número de classes e no mesmo tempo diminui a diferença da altura de colunas seguidas.

Superpostas aos histogramas dos binômios foram desenhadas curvas "normais" ou de GAUSS. Quando $n = 64$, o histograma já acompanha a curva razoavelmente bem.

Em baixo, na figura 1, são dadas as frequências (ordenadas) da curva de GAUSS com classes de 0,25 do erro standard. Se nós compararmos estes valores com as frequências do binômio $(1/2 + 1/2)^{64}$, podemos constatar uma concordância bem acentuada.

A transformação matemática do binômio para a curva da distribuição de GAUSS é dada em detalhe no tratado geral de Estatística de YULES AND KENDALL (1940, pg. 177-180) e chega-se finalmente à seguinte fórmula simples para as ordenadas da distribuição de GAUSS:

$$y = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} \frac{x^2}{p \cdot q \cdot n}} \quad \text{-----} \quad 1-6$$

Podemos ainda substituir alguns dos termos nesta equação. Lembrando que definimos acima x como a distância de m do termo central, podemos também escrever em termos da variação em relação a uma média \bar{v} .

$$x = v - \bar{v}$$

O termo $p \cdot q \cdot n$ é a variância de um binômio, e devemos substituí-lo pela letra σ que caracteriza em geral o erro standard de uma distribuição. Assim chegamos finalmente à fórmula:

$$y = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{v - \bar{v}}{\sigma}\right)^2} = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} D^2} \quad \text{-----} \quad 1-7$$

Esta é a equação que define a curva normal ou de GAUSS, a qual pode ser construída facilmente, calculando-se as ordenadas y para todas as abscissas D pela fórmula 1-7.

A condição principal que permitirá chegar, da equação das frequências do termo binomial (fórmula 1-5) à equação das ordenadas da curva de GAUSS (fórmula 1-7) consiste em supor que o número n atingiu o valor infinito.

Dissemos acima que o número n é o número de pregos do aparelho de GALTON que uma bolinha encontra no seu caminho e que alteram a sua marcha. Estes pregos representam os fatores de acaso que podem afetar um acontecimento qualquer. Não é suficiente porém que o número de pregos seja grande, se nós também não usarmos ao mesmo tempo um número elevado de bolinhas. Se houvesse poucas bolinhas, elas apenas encontrariam alguns pregos no seu caminho, mas para que todos os pregos entrem em jogo, o número de bolinhas deverá

também ser grande. Assim, os pregos representam os fatores de acaso que atuam sobre um acontecimento, neste caso a marcha das bolinhas, e chegamos à conclusão que ambos, os pregos e bolinhas, fatores de acaso e acontecimentos, devem ser muito numerosos, aproximando-se o seu número ao infinito.

O termo "acontecimento" também devemos ainda definir melhor em sua significação geral. Se por exemplo estudarmos a altura de plantas, ou o peso de animais, etc., a altura ou o peso de cada individuo é um tal "acontecimento" que depende de um grande número de causas incontroláveis ou de acaso.

Para poder aplicar a fórmula da distribuição de GAUSS devemos estar certos de que as nossas observações abrangem todas as combinações possíveis dos fatores de acaso e que nenhuma delas ficou sem ser registrada, em consequência de uma limitação do número de observações. Podemos ter esta certeza quando o número de observações é tão grande de modo que não foi feita uma seleção qualquer, dirigida ou não, e que poderia resultar numa omissão de alguns acontecimentos. Para ter certeza que o jogo do número infinitamente grande de agentes de acaso era absolutamente livre e ilimitado e que tudo foi também registrado, o número de observações também deverá ser infinitamente grande.

Como se pode deduzir da fórmula 1-6 e 1-7, o erro standard é o único termo na equação, onde o valor n ainda aparece. Assim, podemos frisar que uma variação de acaso segue a norma estabelecida pela distribuição de GAUSS somente quando, no termo $D = (\sqrt{v} - \bar{v}) : \sigma$ o valor do erro standard for derivado de um número ilimitado ou infinitamente grande de observações.

Em geral, em nossos experimentos, o número de observações ou de variáveis é pequeno, e além disso introduzimos frequentemente durante o processo da análise limitações que restringem a liberdade da variação. Na análise de séries de observações sobre variáveis em número limitado e restringido, isto é, onde não somente o número total, mas também o número de variáveis livres ou o seu grau de liberdade é reduzido, não podemos esperar que a variação siga a distribuição de GAUSS, e as frequências em tais séries não corresponderão às ordenadas desta distribuição. Devemos procurar a definição de outras distribuições do acaso nas quais o jogo do acaso não é mais completamente livre, mas onde há uma limitação da variação e do número de variáveis livres.

2 — A ANÁLISE DO ERRO STANDARD E AS CONCLUSÕES DE LEXIS E DE PEARSON

Pelo que dissemos acima e de acôrdo com a equação 1-6 e 1-7, devemos em primeiro lugar discutir a variação de um erro standard e a sua análise.

Sempre que o valor do erro standard de uma distribuição é calculado de um número limitado de observações, temos que admitir que o valor calculado não é uma constante, mas que êle mesmo é variável no sentido de que determinações repetidas darão valores diversos, todos êles sendo diferentes do valor ideal desconhecido. Assim, é necessário achar uma média da variação do erro standard. Desde ha muito se conhece uma fórmula correspondente ao "erro do erro" e a sua derivação matemática é dada no tratado de Yules and Kendall(1940) pags. 395-401), sendo êle definido pela equação:

$$\sigma[\sigma] = \frac{\sigma}{\sqrt{2n}} \sqrt{\frac{\beta_2 - 1}{2}} \quad \text{----- 2-1}$$

sendo o indice da excessividade ou curtosis:

$$\beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{\sum(v-\bar{v})^4}{(\sum(v-\bar{v})^2)^2}$$

O valor dêste indice é igual a 3 em certos casos, e especialmente no caso da distribuição de Gauss. Sômente nêstes casos temos :

$$\frac{\beta_2 - 1}{2} = \frac{3 - 1}{2} = 1.$$

Desaparecendo a raiz na equação (2-1), podemos escrever:

$$\sigma[\sigma] = \pm \frac{\sigma}{\sqrt{2n}} \quad \text{----- 2-2}$$

para: $\beta_2 = 3$, o que se obtem na distribuição de Gauss e quan-

Como explicámos acima, a variação de uma variável pode ser medida pelo seu desvio relativo $D = (\bar{v} - \bar{v}) : \sigma$ e podemos assim analisar vários desvios isoladamente. Querendo porém resolver, se um conjunto dêles pode ser causado pelo acaso, nós reunimos os seus quadrados, seguindo depois dois caminhos que são apenas algebricamente diferentes. Dividindo a soma pelo quadrado do erro ideal e extraíndo a raiz obtemos o valor que denominamos X :

$$X = \sqrt{\frac{\sum(v - \bar{v})^2}{\sigma_0^2}}$$

ou nós dividimos ainda pelo grau de liberdade correspondente nf , obtendo um desvio relativo :

$$D = \sqrt{\frac{\sum(v - \bar{v})^2}{nf} : \sigma_0^2} = \frac{\sigma_1}{\sigma_0} \text{ ----- } 2-4$$

sendo assim :

$$D = X : \sqrt{nf}$$

PEARSON (1900) determinou a fórmula da variação do termo X e nós chamamos distribuições de Pearson às distribuições correspondentes de

$$D' = \frac{X}{\sqrt{nf}}$$

Na fórmula correspondente aparece o valor nf do grau de liberdade do conjunto de desvios, de modo que para cada valor de nf temos uma distribuição. Existe apenas uma distribuição de GAUSS, porém muitas distribuições de PEARSON.

A descoberta de PEARSON (1900) permitiu depois a STUDENT (1917) achar a solução do caso da variação de um só desvio, quando o erro standard correspondente é apenas baseado em poucas observações, sendo relativamente pequeno o grau de liberdade, correspondente ao erro.

Finalmente, R. A. FISHER deu a solução final e bem geral do problema da variação de desvios, seja individual ou em conjunto, achando a fórmula geral de variação do acaso de um quociente de dois erros standard, cada um com graus de liberdade n_1 e n_2 diferentes:

$$D = \frac{\sigma_1}{\sigma_2} \quad \frac{n_1 = \dots}{n_2 = \dots}$$

3 — A DEFINIÇÃO MATEMÁTICA DOS QUATRO TIPOS DE DISTRIBUIÇÃO DE ACASO

a) As fórmulas das distribuições

As fórmulas que definem os quatro tipos de distribuição de acaso — GAUSS, PEARSON, STUDENT e FISHER — permitem calcular as ordenadas para as respectivas abcissas, sendo as ordenadas frequências e as abcissas desvios relativos. A maioria dos tratados de estatística não apresenta as fórmulas na sua forma mais geral. Os testes derivados e as tabelas correspondentes contêm os valores de abcissas, que correspondem a diferentes limites de probabilidade, em forma bastante alterada na sua expressão algébrica. Demonstraremos em seguida as relações matemáticas das distribuições, na sua forma mais geral.

Nos desvios relativos, podemos distinguir o grau de liberdade do numerador n_1 e do denominador n_2 e notamos neste ponto a situação seguinte:

Fisher

$$1 < n_1 < \text{infinito}$$

$$1 < n_2 < \text{infinito}$$

Student

$$n_1 = 1$$

$$1 < n_2 < \text{infinito}$$

Pearson

$$1 < n_1 < \text{infinito}$$

$$n_2 = \text{infinito}$$

Gauss

$$n_1 = 1$$

$$n_2 = \text{infinito}$$

Limitamos a definição das distribuições de FISHER aos graus de liberdade maiores do que 1 e menores do que infinito. Quando o grau de liberdade n_1 do numerador é igual a 1 e aquele do denominador qualquer valor maior do que 1 e menor do que infinito, encontramos uma das distribuições de STUDENT. De outro lado, quando os limites do grau de liberdade n_1 são os mesmos como na distribuição de FISHER, e os graus de liberdade n_2 são infinitamente grandes, então trata-se de uma das distribuições de PEARSON. Com outras palavras, se na distribuição de FISHER alterarmos apenas o grau de liberdade do numerador para 1, ficando sem alteração o valor de n_2 , passamos para uma distribuição de STUDENT. E se não mudarmos os limites de n_1 , ficando n_2 infinitamente grande, obtem-se uma das distribuições de PEARSON. Finalmente, quando fizermos concomitantemente ambas as alterações ($n_1 = 1/n_2 = \text{infinito}$) obtemos a distribuição normal ou de GAUSS.

Nas tabelas que são frequentemente usadas para o Z-teste (FISHER), F-teste (SNEDECOR) e χ^2 -teste (BRIEGER), a definição das distribuições é mais geral, incluindo os extremos e dando os valores críticos dos testes não somente das distribuições de FISHER, mas para tôdas as quatro distribuições, FISHER, STUDENT, PEARSON e GAUSS simultaneamente.

Nas definições matemáticas dos quatro tipos de distribuição podemos distinguir duas componentes: uma que depende dos valores dos graus de liberdade apenas e por isso é constante para cada distribuição, e outra que contém o valor do desvio relativo D do qual depende o valor da frequência y.

Distribuições de Fisher

$$\text{Grao de liberdade} \quad y = k_1 \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1}{n_2} \cdot D^2\right)^{-\frac{n_1+n_2}{2}} \quad 3-1$$

$$1 < n_1 < \text{inf}$$

$$1 < n_2 < \text{inf}$$

$$k_1 = \frac{\frac{n_1+n_2-2}{2}!}{\frac{n_1-2}{2}! \cdot \frac{n_2-2}{2}!} \cdot 2 \cdot \left(\frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \quad 3-1$$

Distribuições de Pearson •

Grão de liberdade

$$1 < n_1 < \text{inf}$$

$$n_2 = \text{inf.}$$

$$y = k_2 \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2} D^2}$$

$$k_2 = \frac{\frac{n_1}{2}}{\frac{n_1-2}{2}! 2^{\frac{n_1-1}{2}}} \text{ --- 3-2}$$

Distribuições de Student

Grão de liberdade

$$n_1 = 1$$

$$1 < n_2 < \text{inf}$$

$$y = k_3 \cdot \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}}$$

$$k_3 = \frac{\frac{n_2-1}{2}!}{\frac{n_2-2}{2}! \sqrt{n_2 \pi}} \text{ --- 3-3}$$

Distribuições de Gauss

Grão de liberdade

$$n_1 = 1$$

$$n_2 = \text{inf.}$$

$$y = k_4 \cdot e^{-\frac{1}{2} D^2}$$

$$k_4 = \frac{1}{\sqrt{2\pi}} \text{ --- 3-4}$$

b) Cálculo das ordenadas e áreas

As fórmulas dadas acima permitem o cálculo de ordenadas de curvas com área unidade, e tendo como unidade da abscissa o erro standard.

Precisamos explicar o que estes últimos termos significam e também como calculamos a área debaixo de uma curva, isto é, entre a curva e o eixo da abscissa. Para acharmos a área dividimos primeiro a curva de tal modo que as diferenças entre abscissas seguidas seja igual

$$x_1 - x_2 = x_2 - x_3 = x_3 - x_4$$

O tamanho destes intervalos será igual a uma fração do erro standard que podemos calcular, dividindo o erro standard por m , sendo m um número livremente escolhido. A área agora será determinada para cada abscissa aproximadamente, supondo-se que podemos considerar a área total da curva acima de cada intervalo como retângulo. Assim, temos a seguinte relação:

<u>Abcissa</u>	<u>Área</u>	
x_1	$y_1 \cdot \frac{\sigma}{m}$	Área total $\Sigma \left(y \cdot \frac{\sigma}{m} \right)$
x_2	$y_2 \cdot \frac{\sigma}{m}$	

Para uma distribuição com erro standard e área unidade, esta fórmula se transforma na seguinte:

$$\text{Área unidade} = \Sigma \frac{y}{m} = 1 \quad \text{sendo } \sigma = 1$$

Para construir os gráficos deste trabalho, foram calculadas as ordenadas em intervalos de $m = 1:10$, do erro standard.

c) Forma das curvas

Na figura 2 temos quatro distribuições, tôdas representadas em escala igual: a distribuição de GAUSS, a distribuição de STUDENT com $n_1 = 1 / n_2 = 10$, a distribuição de PEAR-

SON com $n_1 = 10 / n_2 = \text{infinito}$ e a distribuição de FISHER para $n_1 = 10 / n_2 = 10$. No centro do gráfico temos a abcissa zero, separando os valores positivos e negativos para a direita e para esquerda. As ordenadas sempre são positivas.

O gráfico deixa aparecer algumas propriedades especiais das curvas.

Tôdas as distribuições, como representadas, são simétricas em relação ao ponto de origem ou valor zero da abcissa. Esta conclusão não precisa de nenhum comentário para as distribuições de STUDENT e de GAUSS. Porém, a referida simetria é apenas de ordem secundária para as outras distribuições que são descontínuas para o valor zero da abcissa, de modo que atualmente temos duas curvas completas e separadas, uma no lado positivo e outra no lado negativo do eixo horizontal. Cada uma destas curvas isoladas é assimétrica e o seu módulo ou máximo não fica longe da abcissa $D = 1$. É muito importante notar que a área de cada curva isolada, de PEARSON e de FISHER, seja no lado direito ou no esquerdo de zero, corresponde apenas à metade da área das distribuições de STUDENT e GAUSS.

d) Comparação matemática das quatro distribuições

Uma análise das fórmulas e dos gráficos (Fig. 2) resulta em estabelecer quatro principais diferenças:

I) Em primeiro lugar vamos determinar as ordenadas para a abcissa D igual a zero :

$$y [\text{Gauss}] = k_4 \cdot e^{-\frac{1}{2}D^2} = k_4$$

$$y [\text{Student}] = k_3 \cdot \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}} = k_3$$

$$y [\text{Pearson}] = k_2 \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2} D^2} = 0 \text{ (zero)}$$

$$y [\text{Fisher}] = k_1 \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1}{n_2} D^2\right)^{-\frac{n_1+n_2}{2}} = 0 \text{ (zero)}$$

Nos primeiros dois casos, a ordenada y para a abcissa zero torna-se igual à constante da equação, k_4 e k_3 . Mas nos outros dois casos a ordenada é zero para a abcissa zero. No gráfico isso se faz sentir na forma já mencionada. As curvas de GAUSS e de STUDENT são contínuas desde mais infinito (canto extremo direito do gráfico) até menos infinito (canto esquerdo). As outras duas curvas, de PEARSON e de FISHER, são descontínuas. Nos gráficos temos duas curvas, separadas; uma à direita e outra à esquerda do centro.

II) O valor do desvio relativo nas distribuições de GAUSS e STUDENT de um lado e naquelas de PEARSON e FISHER representam coisas bem diferentes.

No primeiro caso dividimos um só desvio pelo seu erro, de modo que o quociente, que serve como abcissa no gráfico, tem um sinal definido, positivo ou negativo, de acordo com o sinal do desvio. O valor da ordenada y de outro lado, é independente do sinal de D, pois nas equações 3-3 e 3-4 apenas aparece o quadrado de D que é sempre positivo. Assim, as curvas nos gráficos são simétricas, com uma metade à direita do valor zero da abcissa para os valores positivos do desvio relativo e outra metade à esquerda para os valores negativos.

Nos casos das distribuições de FISHER e PEARSON, o valor do desvio relativo é um quociente de dois erros, respectivamente de um erro por seu valor ideal. A raiz no numerador, e consequentemente também o quociente D poderá ser tanto positivo como negativo. No cálculo como na ilustração gráfica devemos atribuir ao desvio relativo, ou somente o sinal positivo, ou então o sinal negativo. Obtemos assim, como na figura 1, ou uma curva completa à direita do valor zero da abcissa ou uma curva completa à esquerda.

As curvas de PEARSON e FISHER se encontram somente num ou noutro lado do ponto D igual a zero.

III) Em terceiro lugar vamos resolver a questão da posição do máximo das curvas.

Nas distribuições de GAUSS e STUDENT temos que dividir a constante k_4 ou k_3 por um termo exponencial que atinge o seu valor mínimo para D igual a zero, e que cresce para os valores maiores ou menores do que zero. A ordenada y atingirá assim o máximo quando o denominador do quociente é um mínimo, isto é, para a abcissa D igual a zero. O valor da ordenada y para esta abcissa já determinámos acima (sob I) de modo que podemos resumir, dizendo que as curvas de GAUSS e de STUDENT atingem o seu máximo com $y = k_3$ e $y = k_4$, quando D igual a zero.

Nas distribuições de PEARSON e de FISHER a situação é mais complicada, pois as constantes k_2 e k_1 respectivamente devem ser multiplicadas por um termo exponencial de D e divididas por outro termo exponencial de uma função de D . Quando D é pequeno, o primeiro termo cresce mais depressa do que o outro e assim y também cresce, mas depois e na medida que D aproxima-se ao valor 1, o segundo termo cresce mais depressa, e o valor de y decresce conseqüentemente. Explicaremos isso com um exemplo, a distribuição de FISHER com $n_1 = n_2 = 10$.

D	9 log. D	10 log (1+D ²)	Dif. ou log (D ⁹ :(1+D ²) ¹⁰)
0,5	0,28703 — 3	0,96910	0,31163 — 4
0,6	0,00335 — 2	1,33540	0,66795 4
0,7	0,60590 — 2	1,73190	0,87400 4
0,8	0,12772 — 1	2,14840	0,97932 4
0,9	0,58816 — 1	2,57680	0,01136 3
Máximo			
1,0	1,00000	3,01300	0,98700 4
1,1	0,37251	3,44390	0,92861 4
1,2	0,71262	3,87390	0,83872 4

$$\log. k_1 = \log. \left[\frac{9!}{4! 4!} \cdot 2 \right] = 3,10037$$

É evidente neste caso que a diferença entre os logaritmos decimais dos dois termos exponenciais e com isso de y, cresce até $D = 0,9$ caindo depois de novo. Aproximando melhor, achamos o máximo da diferença para $D = 0,904$, sendo o valor do logaritmo desta diferença igual á 0,01151-3. Assim, temos enfim : $y = 0,1294$.

Para a distribuição de FISHER com $n_1 = n_2 = 10$, temos a ordenada máxima de 0,1294 na abcissa $D = 0,904$.

Para a distribuição de PEARSON e FISHER, a ordenada máxima é diferente para cada combinação de graus de liberdade, sendo localizada em abcissas menores do que 1 e aproximando-se ao valor $D = 1$, quando os valores do grau de liberdade se aproximam ao infinito. Uma inspeção dos gráficos na Fig. 2 e nas Figs. 6 a 9 demonstra claramente esta conclusão.

IV) Finalmente, consideremos a questão dos chamados "limites de probabilidade". Em todos os testes da análise estatística aplicamos êste conceito, querendo saber se um desvio relativo está dentro ou fora dos respectivos limites.

A significação dêste termo "limite de probabilidade" podemos facilmente compreender se lembrarmos da definição da área debaixo de uma curva. Esta área é a soma de tôdas as ordenadas que por sua vez são as frequências com as quais são esperados os valores correspondentes da abscissa.

Se erguermos uma ordenada num ponto qualquer da extremidade positiva da curva, dividindo assim a área em duas partes, obteremos a situação seguinte: a soma das frequências de tôdas as abscissas menores daquela que serve como limite e na qual erguemos a ordenada, seja a e a frequência de tôdas as abscissas maiores do que o valor do limite b . Supomos que a abscissa foi escolhida de modo que $a = 95\%$ e $b = 5\%$ da área total, então a soma das frequências de abscissas menores foi 0,95 e a soma das frequências de abscissas maiores apenas 0,05, sendo a área total: $0,95 + 0,05 = 1,00$. O valor correspondente das abscissas chamamos o 0,05 ou 5% **limite unilateral à direita**, ou superior (Figs. 3 e 4).

Se tivéssemos erguido a ordenada na extremidade esquerda da curva, então se trataria do **limite unilateral inferior**

Podemos erguer ordenadas em ambas as extremidades ao mesmo tempo, de modo que a área entre as ordenadas seria

igual a: a e aquela fora de cada uma igual a: $b = \frac{1-a}{2}$, ou no

exemplo dado acima: $a = 0,95$ e $b = \frac{1-a}{2} = 0,025$. Assim te-

mos marcados os **limites bilaterais**, um à esquerda e outro à direita, e a probabilidade total de se obter abscissas maiores do que o limite à direita ou menores do que o limite à esquerda é igual a 0,05 ou 5%.

Uma representação gráfica dêstes três casos: limite unilateral à direita ou superior, limite unilateral à esquerda ou inferior e limites bilaterais para a distribuição de GAUSS e para a distribuição de FISHER com $nf_1 = nf_2 = 10$ encontramos nas figs. 3 e 4. Notamos nestes gráficos duas particularidades:

As abscissas que marcam os limites unilaterais são sempre menores do que aquelas dos limites bilaterais.

Na distribuição de GAUSS (Fig. 3), e o mesmo encontramos nas distribuições de STUDENT, tanto os dois limites unilaterais como também os dois limites bilaterais são de tamanho igual, mas com sinal oposto. Os limites à esquerda ou inferior-

res têm um sinal negativo e os limites à direita ou superiores têm o sinal positivo.

A situação é bastante diferente nas distribuições de PEARSON e de FISHER. (Fig. 4). Aqui não encontramos diferenças de sinais, pois os limites à direita são abcissas maiores do que 1 e os limites à esquerda são frações menores do que 1. Tratando-se de distribuições assimétricas, não existe uma relação matemática simples entre os valores de ambos os limites, exceto em alguns casos especiais que nós iremos mencionar apenas sem entrar numa discussão matemática.

a) Chamamos de distribuições recíprocas aquelas nas quais os graus de liberdade são trocados, como $n_1 = a / n_2 = b$ e $n_1 = b / n_2 = a$. Nestas distribuições recíprocas, os limites superiores de uma distribuição são os valores recíprocos dos limites inferiores das distribuições correspondentes.

b) Nas distribuições de FISHER com $n_1 = n_2$ os valores dos limites inferiores e superiores também são valores recíprocos.

c) Quando os valores dos graus de liberdade são grandes, as distribuições correspondentes tornam-se praticamente simétricas, e os limites superiores e inferiores se tornam equidistantes do valor 1.

Analisando as distribuições de GAUSS ou de STUDENT usamos em geral limites que separam as extremidades em ambos os lados da curva, quando nas distribuições de PEARSON e FISHER consideramos apenas **uma extremidade**, isto é, o limite unilateral superior. Isto é de um modo geral bem justificável.

Quando estudamos um desvio simples D igual a um desvio dividido pelo seu erro standard, admitimos que tanto valores positivos e negativos **devem** acontecer. Uma variação excessiva deve-se fazer sentir nos dois extremos da curva, tanto no lado positivo como no negativo.

Na distribuição de PEARSON, de outro lado, com n_1 entre

um e infinito e $n^2 = \text{infinito}$, comparamos os dois erros σ_1 e σ_0 . O último erro, derivado, seja de uma fórmula teórica ou de um número ilimitado de observações, representa então o erro inevitável do experimento em aprêço.

Assim, um valor do quociente dos erros maior do que 1, indica que a variação na amostra, que deu o valor do primeiro erro, foi excessivo, em consequência de fontes de variação adicionais.

Um valor menor do que 1 não podemos esperar, nem precisamos tomar em consideração, pois se de fato o valor σ_0 é a medida da variação proveniente do acaso σ qual não pode ser excluída do experimento, um valor σ_1 , significativamente menor do que σ_0 não podia ser encontrado.

Seria apenas possível que a variação fôsse aumentada pela ação de fatores especiais, e assim justifica-se plenamente considerarmos somente uma extremidade da curva: aquela que corresponde a abcissas maiores do que 1. As tábuas até hoje publicadas se referem apenas aos limites de probabilidade unilaterais e superiores que indicam a improbabilidade de se obter valores da abcissa maiores do que aquêles da tábua.

Na distribuição de FISHER, a situação é muitas vezes a mesma, quando a natureza dos dados nos permite indicar qual o erro que deve ser considerado como experimental e inevitável ou de acaso. Quando o outro erro for estatisticamente menor do que o erro experimental, isto indicaria uma de duas coisas: ou o erro experimental não foi determinado com exatidão suficiente e a sua estimativa era errada e grande demais, ou na parte do experimento ao qual se refere o erro menor, condições especiais causaram uma limitação sensível da variação.

Na genética, comparando a variação em famílias resultantes de um cruzamento, nas gerações F1 e F2, etc., temos uma variação que depende da variação fenotípica e outra causada pela segregação genotípica ou mendeliana. Esta última deve ter um máximo em F2 e um mínimo em F1, atingindo um grau diferente nas famílias de F3, F4, etc.. Assim, teremos uma variação dos erros standard, por família, em dois sentidos, como indicação da maior ou menor segregação mendeliana.

Em geral, quando comparamos a variação num número grande de amostras com um conjunto de muitos erros, devemos tomar em consideração ambos os limites: o limite superior, que nos indica quais os valores máximos de abcissas, maiores do que 1, como também os limites inferiores para as abcissas menores do que 1.

Os limites em tôdas as tábuas para as distribuições de FISHER até hoje publicados se referem sempre a uma extremidade apenas. Se precisamos tomar em consideração ambas as extremidades simultaneamente, não devemos esquecer que então os limites indicados nestas tábuas para o valor de 5%, 1% ou 0,1% de probabilidade (limites unilaterais), correspondem a 10%, 2% e 0,2% de probabilidade (limites bilaterais).

Considerando que em muitos casos é necessário empregar os limites unilaterais das distribuições de STUDENT e os limites bi-laterais das distribuições de PEARSON e FISHER, preparei novas táboas que já estão em impressão. Nessa publicação será explicado com mais detalhe o problema geral dos limites. (BRIEGER, 1945).

4 — A TRANSFORMAÇÃO MATEMÁTICA DAS 4 DISTRIBUIÇÕES

a) Noções matemáticas

Nestas transformações somos obrigados a recorrer a algumas noções gerais do cálculo matemático, as quais podemos resumir da maneira seguinte:

Nas fórmulas aparecem sempre termos fatoriais da forma geral $(n-2) / 2!$ e $(n-1) / 2!$. Sendo os termos fatoriais uma forma de expressão da função gama, eles têm os seguintes valores em dependência do valor n :

Quando $(n-2) : 2$ ou $(n-1) : 2$ são números positivos inteiros:

$$\frac{n-2}{2}! = 1 \cdot 2 \cdot 3 \cdots \frac{n-2}{2} \qquad \frac{n-1}{2}! = 1 \cdot 2 \cdot 3 \cdots \frac{n-1}{2}$$

Quando $(n-2) : 2$ ou $(n-1) : 2$ tem como valores: 0,5 ou 1,5 ou 2,5...:

$$\frac{n-2}{2}! = 0,5 \cdot 1,5 \cdots \frac{n-2}{2} \cdot \sqrt{\pi} \qquad \frac{n-1}{2}! = 0,5 \cdot 1,5 \cdots \frac{n-1}{2} \cdot \sqrt{\pi}$$

Quando $(n-2) : 2$ ou $(n-1) : 2$ são iguais a 1 ou 0:(zero)

$$0! = 1 \qquad 1! = 1$$

Quando $(n-2) : 2$ e' igual a -0,5:

$$(-0,5)! = \sqrt{\pi}$$

Finalmente devemos considerar a possibilidade de n se aproximar ao valor infinito, e neste caso podemos aplicar a fórmula de STIRLING:

$$\lim_{n \rightarrow \infty} n! = \sqrt{2\pi} \cdot n^{n+\frac{1}{2}} \cdot e^{-n}$$

Outra fórmula de bastante importância, derivada da definição da base e dos logaritmos naturais, é:

$$\lim_{m \rightarrow \infty} \left(1 + \frac{a}{m}\right) = e^{\frac{a}{m}}$$

b) Transformação geral da parte exponencial

Como uma primeira informação daremos a transformação da parte exponencial das equações 3-1 a 3-4.

Quando o grau de liberdade n_2 cresce até infinito temos a seguinte transformação:

$$\begin{aligned} \lim_{n_2 \rightarrow \infty} Y [Fisher] &= k_1 \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1}{n_2} D^2\right)^{-\frac{n_1+n_2}{2}} \\ &= k \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{n_2} D^2} \cdot \frac{n_1+n_2}{2} \end{aligned}$$

igualando n_2 a $(n_2 + n_1)$ quando n_2 infinitamente grande, de modo que $(n_2 + n_1) : n_2 = 1$, temos:

$$\lim_{n_2 \rightarrow \infty} Y [Fisher] = k_2 \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2} D^2} = Y [Pearson] \quad \text{--- 4-1}$$

Do outro lado, partindo-se da distribuição de PEARSON (3-2) e se nós atribuirmos a n_1 o valor de unidade, teremos:

$$\lim_{n_1=1} Y [Pearson] = k_2 \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2} D^2} = k_4 \cdot e^{-\frac{1}{2} D^2} = Y [Gauss] \quad \text{--- 4-2}$$

Atribuindo na equação da distribuição de FISHER (3-1) ao grau de liberdade n_1 o valor de um, a equação se transforma na equação da distribuição de STUDENT (3-3):

$$\lim_{n_1 \rightarrow 1} Y[\text{fisher}] = k_1 \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1 D^2}{n_2}\right)^{-\frac{n_1+n_2}{2}} = k_3 \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}} = Y[\text{Student}] \quad \text{--- 4-3}$$

Finalmente, se na distribuição de STUDENT o valor de n_2 se aproxima ao infinito, obteremos:

$$\lim_{n_2 \rightarrow \infty} Y[\text{Student}] = k_3 \cdot \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}} = k_4 \cdot e^{-\frac{1}{n_2} D^2 \cdot \frac{n_2+1}{2}} = k_4 \cdot e^{-\frac{1}{2} D^2} = Y[\text{Gauss}]$$

pois podemos igualar n_2 e (n_2+1)

Neste capítulo explicamos apenas a transformação matemática da parte das equações que contém as funções do desvio relativo D . Para completar porém o nosso trabalho será ainda necessário demonstrar que os coeficientes k_1 , k_2 , k_3 e k_4 podem facilmente ser transformados, o que nós faremos nos capítulos seguintes.

c) Transformação da equação de FISHER para a de PEARSON

Precisamos demonstrar que a equação (3-1) passa para a equação (3-2) quando n_2 se aproxima ao valor infinito.

Separaremos em primeiro lugar dois grupos de termos na constante k_1 : um primeiro contendo apenas n_1 e que não sofrerá alterações, sendo independente de n_2 e D ; um segundo que se altera com a aproximação do valor de n_2 ao infinito, sendo também independente de D :

$$k_1 = \frac{n_1^{\frac{n_1}{2}}}{\frac{n_1-2}{2}!} \cdot \frac{\frac{n_1+n_2-2}{2}! 2}{\frac{n_2-2}{2}! \frac{n_1}{2}}$$

Na transformação do segundo termo aplicamos os princípios da matemática acima mencionados. Escrevamos ainda para simplificar:

$$a = n_2 - 2$$

$$k' = \frac{\frac{n_1+a}{2}! 2}{\frac{a}{2}! \frac{n_1}{2}}$$

Aplicando agora a fórmula de STIRLING, temos quando n_2 , e também $a = n_2 - 2$ se aproximam ao infinito:

$$\lim_{a \rightarrow \text{inf.}} \frac{n_1 + a}{2}! = \sqrt{2\pi} \cdot \left(\frac{n_1 + a}{2}\right)^{\frac{n_1 + a + 1}{2}} \cdot e^{-\frac{n_1 + a}{2}}$$

$$\lim_{a \rightarrow \text{inf.}} \frac{a}{2}! = \sqrt{2\pi} \cdot \left(\frac{a}{2}\right)^{\frac{a + 1}{2}} \cdot e^{-\frac{a}{2}}$$

$$\lim_{a \rightarrow \text{inf.}} k' = \frac{\sqrt{2\pi} (n_1 + a)^{\frac{n_1 + a + 1}{2}} \cdot 2^{\frac{a + 1}{2}} \cdot e^{\frac{a}{2}} \cdot 2^{\frac{n_1}{2}}}{\sqrt{2\pi} \cdot 2^{\frac{n_1 + a + 1}{2}} \cdot a^{\frac{a + 1}{2}} \cdot e^{\frac{n_1 + a}{2}} \cdot n_2^{\frac{n_1}{2}}}$$

$$= \frac{(n_1 + a)^{\frac{n_1}{2}} \cdot (n_1 + a)^{\frac{a + 1}{2}} \cdot 2^{\frac{a + 1}{2}} \cdot 2 \cdot e^{\frac{a}{2}}}{n_2^{\frac{n_1}{2}} \cdot a^{\frac{a + 1}{2}} \cdot 2^{\frac{n_1 + a + 1}{2}} \cdot e^{\frac{n_1 + a}{2}}}$$

$$= \left(\frac{n_1 + a}{n_2}\right)^{\frac{n_1}{2}} \cdot \left(\frac{n_1 + a}{a}\right)^{\frac{a + 1}{2}} \cdot 2^{\frac{n_1 - 1}{2}} \cdot e^{-\frac{n_1}{2}}$$

$$= \left(\frac{n_1 + a}{n_2 + \frac{a}{n_2}}\right)^{\frac{n_1}{2}} \cdot \left(\frac{n_1 + 1}{a} + 1\right)^{\frac{a + 1}{2}} \cdot 2^{\frac{n_1 - 2}{2}} \cdot e^{-\frac{n_1}{2}}$$

Igualemos $n_2 - 1$, $n_2 - 2 = a$ e n_2 , o que é justificado quando este último valor se aproxima ao infinito, e igualemos também $(n_2 + n_1)$ a n_2 que é permissível quando n_2 é um número muito grande e n_1 um número razoavelmente pequeno, e obteremos:

$$\lim (k') = \left(1 + \frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \cdot \left(1 + \frac{n_1}{n_2}\right)^{\frac{n_2}{2}} \cdot 2 \cdot 2^{-\frac{n_1-2}{2}} \cdot e^{-\frac{n_1}{2}}$$

$$n_2 = \text{inf}$$

$$= \left(1 + \frac{n_1}{n_2}\right)^{\frac{n_1+n_2}{2}} \cdot e^{-\frac{n_1}{2}} \cdot 2^{-\frac{n_1-2}{2}}$$

$$= e^{\frac{n_1}{n_2} \cdot \frac{n_1+n_2}{2}} \cdot e^{-\frac{n_1}{2}} \cdot 2^{-\frac{n_1-2}{2}}$$

$$\approx e^{\frac{n_1}{2}} \cdot e^{-\frac{n_1}{2}} \cdot 2^{-\frac{n_1-2}{2}} = 2^{-\frac{n_1-2}{2}}$$

Voltando agora para a equação inicial de k_1 , teremos:

$$\lim_{n_2 = \text{inf}} (k_1) = \frac{n_1^{\frac{n_1}{2}}}{\left(\frac{n_1-2}{2}\right)!} \cdot \frac{1}{2^{\frac{n_1-2}{2}}} = k_2 [\text{Pearson}] \quad \text{--- --- --- 4-5}$$

Conseguimos assim demonstrar que a constante de FISHER passa para aquela de PEARSON quando n_2 se aproxima a um valor infinitamente grande. Já mostrámos a transformação do resto da equação da distribuição de FISHER (4-1) e podemos escrever finalmente, combinando as equações 4-1 e 4-5:

$$\lim_{n_1=1} Y [Fisher] = \frac{n_1+n_2-2!}{n_1-2! \cdot n_2-2!} \cdot 2 \cdot \left(\frac{n_1}{n_2}\right)^{n_1-1} \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1}{n_2} D\right)^{-\frac{n_1+n_2}{2}}$$

$$= \frac{n_1}{n_1-2! \cdot n_2-2!} \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2} D}$$

$$\lim_{n_2 \rightarrow \infty} Y [Fischer] = Y [Pearson] \quad \text{-----} \quad 4..6$$

d) Transformação da distribuição de PEARSON para a de GAUSS

Temos agora a tarefa de demonstrar que a fórmula 3-2 passa para a fórmula 3-4 quando $n_1 = 1$.

A constante k_2 da fórmula 3-2 tomará a seguinte forma quando substituirmos n_1 pelo valor um:

Conseguimos então demonstrar que a equação da distribuição de FISHER (equação 3-1) passa para a equação da distribuição de PEARSON (equação 3-2) quando n_2 se aproxima a um valor infinitamente grande.

$$k_2 = \frac{n_2}{\frac{n_1-2}{2}! 2^{\frac{n_1-2}{2}}} = \frac{1}{\left(-\frac{1}{2}\right)! 2^{-\frac{1}{2}}} = \sqrt{\frac{2}{\pi}} \quad \text{--- --- --- 4-7}$$

Substituindo também n_1 por um nos outros termos, (equação 4-2), temos:

$$\lim y = \sqrt{\frac{2}{\pi}} \cdot e^{-\frac{1}{2}D^2}$$

Agora devemos lembrar um ponto já discutido acima sob 4b, isto é, que uma curva de PEARSON corresponde apenas à metade de uma distribuição de GAUSS. Para evitar que a área da última ficasse o dobro da primeira, devemos dividir as ordenadas por 2:

$$\lim_{n_1=1} y [Pearson] = \frac{1}{2} \sqrt{\frac{2}{\pi}} \cdot e^{-\frac{1}{2}D^2} = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}D^2} = y [Gauss] \quad \text{--- --- --- 4-8}$$

e) Transformação da equação de FISHER para a de STUDENT

Mostraremos que basta substituir na equação da distribuição de FISHER n_1 por 1, para se obter a fórmula de STUDENT:

$$\begin{aligned} \lim_{n_1=1} y [Fisher] &= \frac{\frac{n_1+n_2-2}{2}!}{\frac{n_1-2}{2}! \frac{n_2-2}{2}!} \cdot 2 \cdot \left(\frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \cdot \left(1 + \frac{n_1}{n_2} D^2\right)^{-\frac{n_1+n_2}{2}} \\ &= \frac{\frac{n_2-1}{2}!}{\left(-\frac{1}{2}\right)! \frac{n_2-2}{2}!} \cdot 2 \cdot \left(\frac{1}{n_2}\right)^{\frac{1}{2}} \cdot \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}} \\ &= \frac{\frac{n_2-1}{2}!}{\frac{n_2-2}{2}! \sqrt{n_2 \pi}} \cdot 2 \cdot \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}} \end{aligned}$$

Considerando que a área da curva de FISHER corresponde apenas à metade da área das curvas de STUDENT como explicado acima (4b) devemos, para igualar as áreas, dividir as ordenadas por dois:

$$\lim_{n_2=1} y[\text{Fisher}] = \frac{\frac{n_2-1}{2}!}{\frac{n_2-2}{2}! \sqrt{n_2 \pi}} \cdot \left(1 + \frac{1}{n_2}\right)^{-\frac{n_2+1}{2}} = y[\text{Student}]$$

----- 4-9

f) Transformação da equação de STUDENT para a de GAUSS quando n_2 aproxima-se ao infinitamente grande

Fazendo a transformação da constante k_3 :

$$\lim_{n_2 = \text{inf}} \left(\frac{n_2-1}{2}\right)! = \sqrt{2\pi} \cdot \left(\frac{n_2-1}{2}\right)^{\frac{n_2}{2}} e^{-\frac{n_2-1}{2}}$$

$$\lim_{n_2 = \text{inf}} \left(\frac{n_2-2}{2}\right)! = \sqrt{2\pi} \cdot \left(\frac{n_2-2}{2}\right)^{\frac{n_2-1}{2}} e^{-\frac{n_2-2}{2}}$$

$$\lim_{n_2 = \text{inf}} k_3 = \frac{\sqrt{2\pi} \cdot (n_2-1)^{\frac{n_2}{2}} \cdot 2^{-\frac{n_2-2}{2}} \cdot e^{\frac{n_2-2}{2}}}{\sqrt{2\pi} \cdot 2^{\frac{n_2}{2}} \cdot (n_2-2)^{\frac{n_2-1}{2}} \cdot e^{\frac{n_2-1}{2}}} \cdot \frac{1}{\sqrt{n_2 \pi}}$$

$$= \frac{(n_2-1)^{\frac{n_2}{2}} \cdot 1 \cdot 1}{(n_2-2)^{\frac{n_2-1}{2}} \cdot 2^{\frac{1}{2}} \cdot e^{\frac{1}{2}} \cdot n_2^{\frac{1}{2}}} \cdot \frac{1}{\sqrt{\pi}}$$

$$= \frac{\left(1 - \frac{1}{n_2}\right)^{\frac{n_2}{2}} \cdot n_2^{\frac{n_2}{2}}}{\left(1 - \frac{2}{n_2}\right)^{\frac{n_2-1}{2}} \cdot n_2^{\frac{n_2-1}{2}} \cdot n_2^{\frac{1}{2}} \cdot e^{\frac{1}{2}}} \cdot \frac{1}{\sqrt{2\pi}}$$

$$= \frac{e^{\frac{1}{2}}}{e^{\frac{n_2-1}{n_2}} \cdot e^{\frac{1}{2}}} \cdot \frac{1}{\sqrt{2\pi}}$$

Sendo, porém, t muito pequeno em relação a n_2 , podemos igualar n_2 a (n_2-1) e obtermos finalmente:

$$\lim_{n_2 = \infty} k_3 = \frac{1}{\sqrt{2\pi}} = k_4 \quad \text{--- 4-10}$$

O segundo termo da equação de STUDENT se transforma simultaneamente, como já demonstrámos acima (4-3), e assim chegamos ao fim da transformação:

$$\lim_{n_2 = \infty} y [\text{Student}] = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}D^2} = y [\text{Gauss}] \quad \text{--- 4-11}$$

5 — A FORMA DAS DISTRIBUIÇÕES EM DEPENDÊNCIA DO GRAU DE LIBERDADE

a) **As distribuições de STUDENT** — A Fig. 5 contém as curvas correspondentes aos valores de $n_2 = 5, 10, 20$ e 30 , sendo n_1 sempre igual a 1. Nota-se que as curvas se tornam cada vez mais altas e mais curtas em ambas as extremidades.

Quando n_2 vai além de 30, podemos usar os limites da distribuição de GAUSS como uma aproximação suficiente. A razão para esta recomendação é bem evidente. A diferença entre as curvas de $n_2 = 20$ e $n_2 = 30$ são da mesma ordem como aquela entre $n_2 = 30$ e $n_2 = \infty$, isto é, a curva de GAUSS.

b) **As distribuições de PEARSON** — A Fig. 6 representa as curvas que correspondem a $n_1 = 2, 5, 10, 20$, sendo n_2 sempre infinito. As curvas são altamente assimétricas quando n_1 é baixo. Quando o grau de liberdade cresce, a curva torna-se mais alta, mais estreita e menos assimétrica.

As curvas devem se aproximar, como explicado no capítulo 2 (Fórmula 2-3) acima, à curva modificada de GAUSS com média $D = 1$, e um erro standard que em vez de unidade, será dividido pela raiz quadrada de $2 n_1$.

Para explicar melhor daremos a fórmula de distribuição de GAUSS com uma pequena modificação: para acentuar que se trata de uma distribuição com erro standard unidade, es-

creveremos, em vez de D , $D:1$, logo, substituindo, teremos a fórmula de uma distribuição modificada de GAUSS:

$$\lim_{n_1 > 30} y \left[\text{Pearson} \right] = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{D^2}{\sqrt{2n_1}} \right)^2} = \frac{1}{\sqrt{2\pi}} \cdot e^{-n_1 \omega^2} \quad \text{---5---1}$$

A Fig. 5, em baixo, ilustra a diferença da curva de PEARSON e da curva modificada de GAUSS, para $n_1 = 20$ e $n_2 = 30$. Especialmente as últimas duas curvas tornam-se já muito semelhantes, de modo que parece justificado, para graus de liberdade superiores a 30, usar em vez da distribuição de PEARSON, aquela de GAUSS.

Podemos determinar os limites de probabilidade pela equação:

$$\left. \begin{array}{l} \text{limites de probabilidades das} \\ \text{distribuições de Pearson para} \\ \text{valores de } n_1 \text{ maiores que 30} \end{array} \right\} = 1 \pm \text{limites (Gauss)} : \sqrt{2n_1} \quad \text{---5---2}$$

Fisher dá à esta relação uma forma matemática diferente. Ele escreve:

$$\sqrt{2X^2} - \sqrt{2n_1 - 1} = \text{limites (Gauss)}$$

Desprezando 1 em relação a $2n_1$ e dividindo por $\sqrt{2n_1}$ teremos:

$$\sqrt{\frac{X^2}{n_1}} - 1 = \text{limite (Gauss)} : \sqrt{2n_1}$$

$$\sqrt{\frac{X^2}{n_1}} = D = 1 \pm \text{limite (Gauss)} \cdot \sqrt{2n_1}$$

o que é matematicamente o mesmo do que a equação dada acima (5-2).

Os limites de GAUSS que nós devemos usar são limites uni-

laterais, pois referem-se apenas a uma extremidade da distribuição, de acôrdo com o que explicámos acima, e são para 5% — 1,64; 1% — 2,33; 1 0/00 — 3,09.

c) **Distribuição de FISHER** — Para melhor esclarecer a situação, trataremos três casos separadamente.

I — Na Fig. 7 encontramos uma série de curvas para as seguintes combinações: $n_1 = n_2 = 2; 5; 10; 20; 30; 50$.

As curvas que correspondem a graus de liberdade baixos são muito assimétricas, mas quando este valor sobe a 30 ou além, a assimetria já parece muito pouco acentuada (Fig. 8), e evidentemente podemos esperar uma aproximação a uma distribuição modificada de GAUSS. Simultaneamente altera-se a curtosis, e as distribuições tornam-se cada vez menos chatas com o crescimento do grau de liberdade. Explicamos acima (pag. 328) que o índice de curtosis ou excessividade β_2 atinge neste casos o seu valor 3, típico para a distribuição de GAUSS.

Para definir matematicamente a natureza da distribuição modificada de GAUSS, devemos fazer as seguintes considerações:

Supomos que dois erros σ_1 e σ_2 , que nós queremos comparar, fôsem desvios de acaso de um valor σ_0 . Assim, e para valores de $n_1 = n_2$ bastante elevados, podemos determinar o êrro correnpondente à variação dos erros pela fórmula :

$$\sigma [\sigma] = \frac{\sigma_0}{\sqrt{2n}}$$

O êrro correspondente ao quociente dos dois erros pode ser determinado aplicando a fórmula geral do êrro de um quociente de duas variáveis a e b, independentes, com médias \bar{a} e \bar{b} e com erros σ_a e σ_b , do seguinte modo :

$$\sigma \left[\frac{a}{b} \right] = \pm \frac{\bar{a}}{b} \sqrt{\left(\frac{\sigma_a}{\bar{a}} \right)^2 + \left(\frac{\sigma_b}{\bar{b}} \right)^2} \quad \text{-----} \quad \text{5-3}$$

$$\sigma \left[\frac{\sigma_1}{\sigma_2} \right] = \pm \frac{\sigma_2}{\sigma_0} \sqrt{\left(\frac{\sigma_0}{\sigma_0 \sqrt{2n}} \right)^2 + \left(\frac{\sigma_0}{\sigma_0 \sqrt{2n}} \right)^2} = \pm \frac{1}{\sqrt{n}} \quad \text{---} \quad \text{5-4}$$

Assim podemos definir as ordenadas da distribuição modificada de GAUSS do quociente dos dois erros quando $n_1 = n_2 = n$ como sendo igual a:

$$\lim y [fisher] = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{D^2}{\sqrt{n}}\right)^2} = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{D^2}{2n}} \quad \text{---5---5}$$

Os limites da probabilidade podem assim ser calculados quando $n_1 = n_2 = n$, e sendo n um valor grande pela fórmula simples:

$$\frac{G_1}{G_2} = 1 \pm (\text{limites de Gauss}) : \sqrt{n} \quad \text{---5---6}$$

Os limites convencionais da distribuição de GAUSS são:

	5%	1%	10/100
para uma extremidade	1,64	2,33	3,09
para ambas as extremidades	1,96	2,53	3,29

Na Fig. 8 estão reproduzidas as duas curvas de FISHER para $n_1 = n_2 = 30$ e $n_1 = n_2 = 50$, em comparação com as respectivas curvas modificadas de GAUSS. Evidentemente a aproximação não é suficiente para $n_1 = n_2 = 30$, e apenas moderada para $n_1 = n_2 = 50$. Como melhor aproximação podemos recomendar aplicar a fórmula dada acima (5-6) quando $n_1 = n_2$ é igual ou superior a 100.

II — Em segundo lugar estudaremos casos onde os valores de n_1 e n_2 são diferentes.

A parte superior da figura (Fig. 9) contém as curvas que têm todas o mesmo valor $n_1 = 10$, nas seguintes combinações:

n_1	10	10	10	10	10	10	10	10
n_2	1	2	5	10	20	50	100	infinito

A primeira destas distribuições é evidentemente uma curva recíproca de STUDENT para $n_1 = 10/n_2 = 1$ em vez de $n_1 = 1 / n_2 = 10$. A fórmula desta nova distribuição podemos facilmente obter por uma transformação semelhante àquela do capítulo 5e.

Transformemos em primeiro lugar o último termo da equação 3-1 para $n_2 = 1$.

$$\frac{D^{n_1-1}}{(1+n_1 D^2)^{\frac{n_1+1}{2}}} = \frac{D^{n_1-1} \cdot (n_1 D^2)^{-\frac{n_1+1}{2}}}{\left(1 + \frac{1}{n_1 D^2}\right)^{\frac{n_1+1}{2}}} = \frac{1}{n_1} \cdot D^{-2} \cdot \left(1 + \frac{1}{n_1 D^2}\right)^{-\frac{n_1+1}{2}}$$

Transformando agora também a constante k_1 , e reagrupando os termos da área da curva além do valor D igual a zero, obteremos para a distribuição recíproca de STUDENT, com:

$$\begin{aligned} \lim_{n_2=1} Y [fisher] &= \frac{\frac{n_1-1}{2}!}{\frac{n_1-2}{2}! \sqrt{\pi}} \cdot \frac{n_1^{\frac{n_1}{2}}}{n_1} \cdot D^{-2} \cdot \left(1 + \frac{1}{n_1 D^2}\right)^{-\frac{n_1+1}{2}} \\ &= \frac{\frac{n_1-1}{2}!}{\frac{n_1-2}{2}! \sqrt{n_1 \pi}} \cdot D^{-2} \cdot \left(1 + \frac{1}{n_1 D^2}\right)^{-\frac{n_1+1}{2}} \end{aligned}$$

De outro lado, a última curva da Fig. 9 (parte superior), com $n_1 = 10/n_2 = \text{infinito}$, é uma distribuição típica de PEARSON.

III — Na Fig. 9, parte inferior, temos representadas as curvas das distribuições recíprocas, sendo agora $n_2 = 10$:

n1		1		2		5		10		20		50		100		infinito
n2		10		10		10		10		10		10		10		10

A primeira curva é evidentemente uma curva de STUDENT, mas para manter as curvas comparáveis, será necessário usar apenas a metade da curva de STUDENT, multiplicando as ordenadas por 2 para ter-se uma área de unidade.

O outro extremo podemos considerar como a curva "recíproca" da distribuição de PEARSON sendo $n_1 = \text{infinito}$, em vez de n_2 . As ordenadas de uma tal curva podem ser determinadas facilmente, depois de uma transformação conveniente da fórmula de FISHER:

$$Y [fisher] = \frac{1}{\frac{n_2-2}{2}!} \cdot \frac{\frac{n_1+n_2-2}{2}! \cdot 2}{\frac{n_1-2}{2}!} \cdot \left(\frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1}{n_2} D^2\right)^{-\frac{n_1+n_2}{2}}$$

O primeiro termo é independente de n_1 . O segundo termo, que podemos chamar K'' , podemos transformar em analogia ao termo k' da equação 4-5.

$$\lim_{n_1 \rightarrow \infty} \left[\frac{\frac{n_1+n_2-2}{2} ! 2}{\frac{n_1-2}{2} ! n_1^{\frac{n_2}{2}}} \right] = \lim k'' \cdot (n_1)^{\frac{n_2}{2}} = \frac{n_1^{\frac{n_2}{2}}}{2^{\frac{n_2-2}{2}}}$$

Temos ainda que tomar em consideração os termos restantes:

$$\begin{aligned} \left(\frac{n_1}{n_2}\right)^{\frac{n_2}{2}} \cdot \frac{D^{n_1-1}}{\left(1 + \frac{n_1}{n_2} D\right)^{\frac{n_1+n_2}{2}}} &= \left(\frac{n_1}{n_2}\right)^{\frac{n_2}{2}} \cdot \frac{D^{n_1-1} \left(\frac{n_2}{n_1 D}\right)^{\frac{n_1+n_2}{2}}}{\left(\frac{n_2}{n_1 D} + 1\right)^{\frac{n_1+n_2}{2}}} \\ &= \left(\frac{n_1}{n_2}\right)^{\frac{n_2}{2}} \left(\frac{n_2}{n_1}\right)^{\frac{n_1+n_2}{2}} D^{-(n_1+n_2)} \left(1 + \frac{n_2}{n_1 D}\right)^{-\frac{n_1+n_2}{2}} \\ &= \left(\frac{n_1}{n_2}\right)^{\frac{n_2}{2}} D^{-(n_2+1)} \left(\frac{n_2}{n_1 D}\right)^{-\frac{n_1+n_2}{2}} \end{aligned}$$

O último termo transforma-se ainda em:

$$\lim_{n_1 \rightarrow \infty} \left(1 + \frac{n_2}{n_1 D}\right)^{-\frac{n_1+n_2}{2}} = e^{-\frac{n_2}{2D}} \cdot \frac{n_2}{n_1 D} \cdot \frac{n_1+n_2}{2} = e^{-\frac{n_2}{2D}} \cdot \frac{n_1+n_2}{n_1} = e^{-\frac{n_2}{2D}}$$

pois sendo n_2 pequeno em relação a n_1 , que é igual a infinito, podemos escrever $(n_1 + n_2) : n_1$ igual a 1.

Assim, chegámos finalmente à seguinte transformação do conjunto, para a equação da distribuição recíproca de PEARSON:

$$\begin{aligned} \lim_{n_1 \rightarrow \infty} \gamma [\text{fisher}] &= \frac{1}{\frac{n_2-2}{2}} \cdot \frac{n_1^{\frac{n_2}{2}}}{2^{\frac{n_2-1}{2}}} \cdot \left(\frac{n_2}{n_1}\right)^{\frac{n_2}{2}} \cdot D^{-(n_2+1)} \cdot e^{-\frac{n_2}{2} D^2} \\ &= \frac{n_1^{\frac{n_2}{2}}}{\frac{n_2-2}{2} \cdot 2^{\frac{n_2-1}{2}}} \cdot D^{-(n_2+1)} \cdot e^{-\frac{n_2}{2} D^2} = \gamma [\text{PEARSON recíproca}] \quad \text{5-8} \end{aligned}$$

Mencionamos acima sob 3a que nós podemos considerar as distribuições de STUDENT e de PEARSON como os casos extremos da distribuição de FISHER, e também explicamos sob 3e que a área das curvas de FISHER corresponde à metade da área das de STUDENT. Agora comprovamos estes pontos de forma matemática mais concisa.

Uma inspeção das curvas na Fig. 9 mostra que a diferença entre as curvas para a combinação 10/50 e 10/100 de um lado, e de 10/100 e 10/infinito são da mesma ordem. Daí podemos concluir que quando um dos graus de liberdade é maior do que 100 nós podemos substituir a distribuição de FISHER pela distribuição correspondente de PEARSON para todos os fins da análise estatística.

Explicámos acima que as distribuições de PEARSON por sua vez podem ser substituídas por uma distribuição modificada de GAUSS quando o grau de liberdade n_1 ultrapassa 30. Agora devemos determinar qual a distribuição modificada de GAUSS que corresponde ao limite de distribuição de FISHER quando n_1 difere de n_2 e quando um deles é igual ou maior do que 50 e o outro do que 100. A média destas distribuições é $D=1$, e a unidade da abcissa será um erro standard que podemos derivar, como no caso anterior de $n_1 = n_2$, da fórmula do quociente (5-5), pela transformação :

$$\begin{aligned} \sigma \left(\frac{\sigma_1}{\sigma_2} \right) &= \pm \frac{\sigma_0}{\sigma_0} \sqrt{\left(\frac{\sigma_0}{\sigma_0 \sqrt{2n_1}} \right)^2 + \left(\frac{\sigma_0}{\sigma_0 \sqrt{2n_2}} \right)^2} \\ &= \pm \sqrt{\frac{n_1 + n_2}{2 \cdot n_1 \cdot n_2}} \quad \text{--- --- --- 5-9} \end{aligned}$$

Assim teremos:

$$\lim Y [\text{fisher}] = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} D^2 \cdot \frac{n_1 + n_2}{2 \cdot n_1 \cdot n_2}} \quad \text{--- 5-10}$$

para $n_1 = 50/n_2 = 100$ ou $n_1 = 100/n_2 = 50$

Os limites de probabilidade podemos assim determinar pela equação:

$$\left. \begin{array}{l} \text{Limites de probabilidade das distribuições de Fisher} \\ \text{para } n_1 > 50 ; n_2 > 100 \\ \text{ou } n_1 > 100 ; n_2 > 50 \end{array} \right\} = 1 \pm \text{limites [Gauss]} : \sqrt{\frac{n_1 + n_2}{2n_1 n_2}}$$

--- 5-11

6 — RELAÇÃO ENTRE AS DISTRIBUIÇÕES E OS TESTES DERIVADOS

Discutimos acima as distribuições de acaso empregando para isso fórmulas estritamente comparáveis, com a mesma

área e usando sempre o valor do desvio relativo D como a abcissa. Vamos mostrar numa publicação futura e já em preparação, que os testes necessários na análise estatística, podem ser derivados todos destas mesmas fórmulas. Porém, os testes como usados da literatura, se referem geralmente a distribuições com outras unidades como abcissas.

Assim, as distribuições de FISHER são em geral usadas hoje com três abcissas diferentes: o termo original z de FISHER (1923), o termo F de SNEDECOR (1937) e o termo ψ de BRIEGER (1937), sendo estes ligados pela equação:

$$D = \frac{\sigma_1}{\sigma_2} = \psi \text{ [Fisher]} = \sqrt{F} \text{ [Snedecor]} = \psi \text{ [Brieger]} \quad \text{---}$$

----- 6-1

Igualmente, a distribuição de PEARSON é aplicada no X^2 teste com a abcissa X^2 , sendo:

$$D = \frac{\sigma_1}{\sigma_2} = \sqrt{\frac{X^2}{n_1}} \quad \text{ou} \quad X^2 = \frac{\sigma_1^2 \cdot n_1}{\sigma_2^2} = D^2 \cdot n_1 \quad \text{-----}$$

----- 6-2

A aplicação dos testes nestas formas derivadas parece-me ter bastante conveniências.

Em primeiro lugar o técnico, quando não especializado, é em geral incapaz de compreender que todos os testes, que empregam os limites das distribuições de acaso, são intimamente ligados de modo que ele aplica os testes sem compreender do que se trata de fato — uma situação bastante perigosa, especialmente na análise de casos complicados.

Em segundo lugar, o técnico acostuma-se a exigir um grau diferente de exatidão nos testes, pois no t-teste, derivado das distribuições de STUDENT, ele analisa os valores simples dos desvios relativos sendo $D = t$, e no F-teste, de SNEDECOR de acordo com as distribuições de FISHER, ele usa os quadrados dos desvios $D^2 = F$. Algebricamente, a segunda casa decimal no t-teste corresponde à quarta casa decimal no F-teste, sendo por exemplo o quadrado de 0,01 igual a 0,0001.

Na aplicação do X^2 teste a situação é pior ainda e por exemplo aqui temos o seguinte, para $n_1 = 10$:

$D = 9,0$	$X^2 = 81,0$		$D = 1,0$	$X^2 = 10$
$D = 0,9$	$X^2 = 8,1$		$D = 0,1$	$X^2 = 0,1$
$D = 0,09$	$X^2 = 0,081$		$D = 0,01$	$X^2 = 0,001$

Unidades de D dão assim dezenas ou centenas em X^2 , decimais de D dão unidades ou décimos de X^2 , centésimos de D dão centésimos ou milésimos de X^2 .

Para o caso da distribuição de PEARSON vamos mostrar como se deve transformar as fórmulas para manter a área unidade, apesar da alteração da unidade da abcissa. A equação para estas distribuições com abcissa D era a equação (3-2):

$$y [Pearson] = \frac{n_1^{\frac{n_1}{2}}}{\frac{n_1-2}{2}! \cdot 2^{\frac{n_1-2}{2}}} \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2}D^2} \quad \text{--- 6-3}$$

Substituímos agora D por $X \cdot n_1^{\frac{1}{2}}$

$$y [Pearson] = \frac{n_1^{\frac{n_1}{2}}}{\frac{n_1-2}{2}! \cdot 2^{\frac{n_1-2}{2}}} \cdot X^{n_1-1} \cdot n_1^{-\frac{n_1-1}{2}} \cdot e^{-\frac{n_1}{2} \cdot X^2 \cdot \frac{1}{n_1}}$$

$$= \frac{\sqrt{n_1}}{\frac{n_1-2}{2}! \cdot 2^{\frac{n_1-2}{2}}} \cdot X^{n_1-1} \cdot e^{-\frac{1}{2}X^2}$$

Para poder agora passar da abcissa D para X sem alterar a área da curva, temos que dividir também y por $\sqrt{n_1}$:

$$y [X] = \frac{1}{\frac{n_1-2}{2}! \cdot 2^{\frac{n_1-2}{2}}} \cdot X^{n_1-1} \cdot e^{-\frac{1}{2}X^2} \quad \text{--- 6-4}$$

Se agora quisermos ainda passar para a distribuição de X^2 em vez de X , teremos ainda que dividir as ordenadas por X :

$$\begin{aligned}
 y [X^2] &= \frac{y [X]}{X} \cdot \frac{y [Pearson]}{X : \sqrt{\eta_1}} = \frac{1}{\frac{\eta_1 - 2}{2} !} \frac{\eta_1 - 2}{2} \cdot X^{\eta_1 - 2} \cdot e^{-\frac{1}{2} X^2} \\
 &= \frac{1}{\frac{\eta_1 - 2}{2} !} \cdot \left(\frac{X^2}{2} \right)^{\frac{\eta_1 - 2}{2}} \cdot e^{-\frac{X^2}{2}}
 \end{aligned}$$

----- 6-5

As três equações 6-3, 6-4 e 6-5 podem ser encontradas na literatura, porém em geral os autores não explicam claramente que se trata da mesma distribuição básica do ponto de vista matemático, porém com diferentes sistemas de ordenadas e abscissas, sendo a unidade das abscissas o valor $D = \sigma_1 : \sigma_0$, na fórmula 6-3, X na fórmula 6-4 e finalmente X^2 na fórmula 6-5.

7 — CONSIDERAÇÕES FINAIS

Os quatro tipos de distribuições de acaso, discutidos nesta publicação, servem como base da análise estatística. Grande foi o mérito da escola inglesa de KARL PEARSON, STUDENT e R. A. FISHER ter mostrado que a chamada distribuição "normal" é apenas um caso especial, com um nome bastante enganador, pois um caso especial nunca podia ser considerado como norma. Com as descobertas destes três autores muitos casos que antes não podiam ser submetidos a uma análise efi-

ciente, puderam finalmente ser estudados e analisados com exatidão.

Não devemos, porém, esperar que agora todos os problemas encontrarão uma solução fácil para a sua análise. Frequentemente temos que submeter antes os dados a transformações algébricas, para poder submetê-los à análise, e estas transformações não raramente envolvem princípios da teoria estatística. Em outros casos a situação torna-se ainda mais complicada, pois trata-se de combinações de variáveis, cada uma das quais segue um outro tipo de distribuição de acaso. Explicaremos estas complicações com algum detalhe.

a) **Distribuição de FISHER** — Quando analisarmos um desvio relativo

$$D = \sqrt{\frac{\sum (v - \bar{v})^2}{n_1}} : \sigma_v$$

podemos aplicar o princípio da decomposição ou da composição. Supomos que temos uma série de m amostras, com médias $\bar{v}_a, \bar{v}_b, \dots$ e graus de liberdade n_a, n_b, \dots . Então podemos decompor o conjunto $\sum (v - \bar{v})^2$ em partes, devendo ser porém os respectivos graus de liberdade todos iguais:

$$n_a = n_b = \dots$$

$$n_a + n_b + \dots = m \cdot n_a = m \cdot n_b$$

$$\frac{\sigma_a^2 + \sigma_b^2 + \dots}{m} = \frac{1}{m} \left(\frac{\sum (v - \bar{v}_a)^2}{n_a} + \frac{\sum (v - \bar{v}_b)^2}{n_b} + \dots \right)$$

$$= \frac{\sum (v - \bar{v}_a)^2 + \sum (v - \bar{v}_b)^2 + \dots}{n_a + n_b + \dots} = \frac{\sum (v - \bar{v})^2}{m \cdot n_p} = \sigma_v^2$$

----- 7-1

Esta derivação pode também ser lida no sentido inverso, e teremos então em vez de uma composição de um erro total, a sua decomposição.

Todos os erros parciais σ_a , σ_b , se baseiam no mesmo grau de liberdade, na $= nb$, de modo que eles foram determinados com a mesma exatidão e podem ser comparados em conjunto, pelo teste de FISHER.

Porém, se os graus de liberdade individuais forem diferentes, o conjunto não mais pode ser considerado como homogêneo. Para a análise de cada erro parcial teremos outro grau de liberdade e assim outros limites calculados de diferentes distribuições de FISHER. Daí a exigência na decomposição do erro ("analysis of variance") que as partes que são reunidas no erro residual simples deveriam ser baseadas no mesmo grau de liberdade.

Uma outra complicação pode ser causada às vezes quando os diferentes erros parciais σ_a , σ_b não podem ser considerados comparáveis. Por exemplo, em experiências de adubação, pode acontecer que a variação dos canteiros com adubo e sem adubo seja diferente. No primeiro entram os enganos na distribuição, mistura e pesagem dos adubos, e nos outros o efeito de elementos mínimos no solo que terão um efeito mais forte por falta de qualquer compensação biológica. Aqui então um teste estatístico de homogeneidade será necessário para provar se as componentes podiam ser combinadas estatisticamente.

b) Depois discutiremos o X^2 teste. Temos às vezes que comparar desvios de dimensão diferente, isto é, diferenças $(\sqrt{v} - \sqrt{v'})$ nas quais os valores de \sqrt{v} são de dimensões bastante diferentes. Assim, PEARSON, quando introduziu o X^2 teste para comparar frequências observadas com frequências esperadas, se viu na dificuldade de reduzi-las a uma base comparável e que permitisse a sua combinação. Por isso devemos no X^2 teste dividir primeiro cada quadrado da diferença entre uma frequência observada e esperada pela última, para depois somar todos estes quocientes.

$$X^2 = \sum \frac{(f_{ob.} - f_{esp.})^2}{f_{esp.}} \quad \text{--- --- --- --- ---} \quad \text{7-2}$$

Se num caso especial todos os valores de frequência esperada fôssem os mesmos, poderíamos escrever:

$$X^2 = \frac{\sum (f_{ob.} - f_{esp.})^2}{f_{esp.}} = \frac{\sum (v - v_0)^2}{v_0} \quad \text{--- --- --- --- ---} \quad \text{7-3}$$

PEARSON (1900) demonstrou que as raízes destes quocientes são distribuídas de acôrdo com a fórmula 6-5 com área unidade e usando X^2 com uma unidade de abcissa. Se extrairmos a raiz quadrada e dividirmos os valores ainda por $\sqrt{D_1}$ eles ficam distribuídos segundo uma distribuição típica de PEARSON, de acôrdo com a nossa fórmula 3-2, com área unidade e usando o erro standard como unidade da abcissa.

c) Uma situação semelhante encontramos com referência à distribuição de FISHER, quando queremos comparar dois erros σ_1 e σ_2 que se referem a amostras com médias bastante diferentes. Este é um problema que aparece frequentemente na genética quando temos que analisar a segregação mendeliana de caracteres quantitativos. Por exemplo, em F2 e em F3 de um cruzamento, teremos nas diferentes famílias médias bastante diferentes em consequência da segregação na geração anterior e devemos investigar se há ainda em cada família indicaçõs da continuação da segregação mendeliana. Devemos assim comparar os respectivos erros entre si ou com aqueles de famílias de gerações anteriores num teste que usa as distribuições de FISHER. Em analogia ao caso de X^2 do capítulo anterior (b) dividimos os desvios quadrados pelas suas médias, ou os erros pela raiz quadrada da média, de modo que obtemos para o desvio relativo:

$$D = \sqrt{\frac{\sigma_1^2 : \bar{u}_1}{\sigma_2^2 : \bar{u}_2}} = \frac{\sigma_1 \sqrt{\bar{u}_2}}{\sigma_2 \sqrt{\bar{u}_1}} \text{ --- --- --- --- --- } 7-4$$

Tratamos as médias como valores constantes, apenas como indicações de dimensão algébrica desprezando o fato que elas também variam entre certos limites.

Assim parece justificado supor que os quocientes do índice da variação ($V_{\sigma^2} : \bar{v}$) variam de acôrdo com uma distribuição de FISHER (BRIEGER 1943).

d) Na análise de desvios simples devemos também resolver se as médias podem ser consideradas como constantes ou não. Nós não podemos considerar o desvio ($v - \bar{v}$) como uma só variável, pois o erro se refere apenas à variação dos valores v e não a diferença ($v - \bar{v}$). De fato, deveríamos aplicar uma fórmula mais complicada:

$$\frac{v - \bar{v}}{\sigma_{dif}} = \frac{v - \bar{v}}{\sqrt{\sigma^2 + \sigma_v^2}} \text{ --- --- --- --- --- } 7-5$$

Existem duas fórmulas para determinar o valor do erro da média. Se pudermos repetir o experimento em idênticas condições m vezes, cada vez obtendo um outro valor para a média \bar{v} , além da média geral de tôdas as repetições $\bar{\bar{v}}$, o erro das médias parciais poderá ser calculado pela fórmula:

$$\sigma_{\bar{v}} = \pm \sqrt{\frac{\sum(\bar{v} - \bar{\bar{v}})^2}{m-1}} \quad \text{--- --- --- 7-6}$$

Podemos neste caso esperar que os valores das diferentes médias parciais variem em redor da média geral $\bar{\bar{v}}$ e com erro standard de acôrdo com uma distribuição de STUDENT com $(m-1)$ graus de liberdade.

De outro lado, quando há uma só série de variáveis, nós não conhecendo o valor de $\bar{\bar{v}}$, temos que aplicar uma fórmula teórica bastante conhecida:

$$\sigma_{\bar{v}} = \pm \frac{\sigma}{\sqrt{N}} \quad \text{--- --- --- 7-7}$$

A derivação da última fórmula inclui a hipótese que a média pode variar em volta de um valor ideal desconhecido, seguindo a distribuição normal ou de GAUSS. Aparece na fórmula 7-10 não o grau de liberdade, mas o número total (N) de variáveis ao qual se refere a média.

Devemos agora substituir os valores das fórmulas 7-6 e 7-7 respectivamente na fórmula 7-5. No último caso obteremos a equação:

$$\frac{v - \bar{v}}{\sigma \text{ dif.}} = \frac{v - \bar{v}}{\sqrt{\sigma^2 + \frac{\sigma^2}{N}}} = \frac{v - \bar{v}}{\sqrt{1 + \frac{1}{N}}} \quad \text{--- --- --- 7-8}$$

Agora temos que resolver dois problemas, um que é relativamente fácil e representa uma simples questão de álgebra, enquanto o outro envolve problemas da teoria estatística.

Quanto ao primeiro podemos dizer que em geral a inclusão do valor do erro da média muito pouco altera o valor final

Este é especialmente claro quando comparamos os dois valores:

$$\left(\frac{v-\bar{v}}{\sigma}\right) \text{ e } \left(\frac{v-\bar{v}}{\sigma} \cdot \frac{1}{\sqrt{1+\frac{1}{N}}}\right)$$

N	$\sqrt{1+\frac{1}{N}}$	$1:\sqrt{1+\frac{1}{N}}$
5	1,10	0,91
10	1,05	0,95
20	1,02	0,98
50	1,01	0,99

O valor do desvio, calculado do modo comum já para $N=5$ é apenas 1,10 vezes maior do que quando nós incluímos o erro da média.

Para o caso do emprêgo da fórmula 7-6 a situação não pode ser resolvida de uma maneira geral. Devemos em cada caso comparar os valores algébricos de σ e $\sigma\bar{v}$ podendo desprezar o último quando fôr pequeno em relação ao primeiro.

O problema estatístico se resolve diferentemente nos dois casos.

Empregando as fórmulas 7-5 e 7-6 em combinação temos um quociente no qual o numerador tem 1 grau de liberdade e o denominador $(N_2-1) + (m-1)$ graus de liberdade, devendo o quociente seguir assim uma distribuição de STUDENT. Este no mínimo é o processo comumente empregado para calcular o grau de liberdade do erro de uma diferença, que porém é apenas justificado quando os dois erros, no nosso caso σ e $\sigma\bar{v}$, são mais ou menos do mesmo tamanho. Quando desprezamos o valor $\sigma\bar{v}$ não devemos tomar em consideração o seu grau de liberdade.

Quando usamos as fórmulas 7-8 e 7-10 em combinação, a situação é mais complicada, pois agora nós combinamos na fórmula do erro da diferença, um erro standard que segue uma distribuição de STUDENT (σ) e outro ($\sigma\bar{v}$) que segue uma distribuição de GAUSS. A combinação dos graus de liberdade pelo modo comum terá o seguinte valor: $(N-1) + \text{infinito}$; assim, o grau de liberdade do erro mais importante σ prática-

mente fica sem efeito por ser sempre pequeno em relação ao valor infinito.

Sendo $\sigma_{\sqrt{v}}$ pequeno em relação a σ , usamos apenas o valor σ , mas sendo o grau de liberdade n pequeno em relação ao grau de liberdade de $\sigma_{\sqrt{v}}$, desprezamos o primeiro. Chegamos assim a um absurdo: combinamos o valor de um erro com o grau de liberdade do outro!

As dificuldades desaparecem quando nós desprezamos por completo tanto o valor como o grau de liberdade do erro das médias. Se julgarmos que isso não é aconselhável, deve-se aplicar a transformação de BEHRENS, que permite combinar uma distribuição de STUDENT com uma de GAUSS (FISHER and YATES 1943).

e) Uma situação semelhante encontramos quando queremos analisar a diferença entre dois valores a e b , que podem ser simples variáveis ou médias, com erro standard σ_a e σ_b , baseados em n_a e n_b graus de liberdade. Devemos calcular como sempre o desvio relativo:

$$\frac{a-b}{\sigma_{dif}} \text{ ----- } 7-9$$

Limitando-nos a diferenças entre valores a e b , não coordenados e não correlacionados, temos duas fórmulas para calcular o erro da diferença, uma que nós podemos chamar a fórmula clássica, e outra a fórmula de FISHER:

Fórmula clássica:

$$\sigma_{dif} = \pm \sqrt{\sigma_a^2 + \sigma_b^2} \text{ ----- } 7-10$$

Fórmula de Fisher

$$\sigma_{\text{dif}} = \pm \sqrt{\sigma_a^2 + \sigma_b^2} = \pm \sigma_o \sqrt{2} \quad \text{--- --- --- --- --- 7-11}$$

$$\sigma_o = \pm \sqrt{\frac{\sum(a - \bar{a})^2 + \sum(b - \bar{b})^2}{n_a + n_b}}$$

$$= \pm \sqrt{\frac{n_a \sigma_a^2 + n_b \sigma_b^2}{n_a + n_b}}$$

A derivação da primeira fórmula é fácil e não envolve nenhum princípio estatístico complicado. Devemos resolver qual distribuição, que a diferença (a-b) medida por este erro standard, deve seguir. De acôrdo com os números n_a e n_b os respectivos erros σ_a e σ_b correspondem a diferentes distribuições de STUDENT. Se queremos aplicar a teoria da estatística com todo rigor, deveríamos recorrer a uma modificação do teste de BEHREND'S mencionado no capítulo anterior. Como uma aproximação satisfatória podemos usar uma distribuição de STUDENT com $(n_a + n_b)$ graus de liberdade, exceto quando um dos graus de liberdade for muito maior do que o outro. Neste último caso encontramos uma situação idêntica àquela discutida no capítulo anterior (7d).

Na fórmula de FISHER (7-11) para o erro de uma diferença encontramos uma situação diferente. Como se pode ver, trata-se de uma média balanceada entre as duas variâncias. Esta fórmula pode ser aplicada apenas quando estamos certos que os dois erros não são estatisticamente diferentes, porque somente neste caso é justificado reunir as variâncias para obter uma única estimativa, por sua vez baseada em $(n_a + n_b)$ graus de liberdade.

A questão da distribuição não oferece agora nenhuma dificuldade. A diferença (a-b) dividida pela estimativa única do seu erro (7-10) deve seguir uma distribuição de STUDENT com $(n_a + n_b)$ graus de liberdade.

Quando, em vez de variáveis simples a e b, se tratar de médias, as duas fórmulas sofrem apenas uma ligeira alteração, quando podemos calcular erros da média pela fórmula: $\sigma: \sqrt{N}$
 Substituindo teremos:

Fórmula clássica:

$$\sigma[\text{dif } \bar{v}] = \pm \sqrt{\frac{\sigma_a^2}{N_a} + \frac{\sigma_b^2}{N_b}} \quad \text{--- 7-12}$$

Fórmula de Fisher:

$$\sigma[\text{dif } \bar{v}] = \pm \sigma_0 \sqrt{\frac{1}{N_a} + \frac{1}{N_b}} \quad \text{--- 7-13}$$

f) A substituição de uma estimativa do erro standard por outra é geralmente feita na análise estatística chamada "analysis of variance", sem que em geral os autores estejam ao par desta substituição. Usamos como exemplo um caso simples de um teste "entre-dentro" ou "betwen-within".

Supomos que temos m amostras, cada uma com (np + 1) observações ou medidas, e queremos saber se as médias

\bar{v}_1 e \bar{v}_2 são diferentes da média geral \bar{v} .

Poderíamos calcular o erro para cada amostra:

$$\sigma_1 = \pm \sqrt{\frac{\sum(v - \bar{v}_j)^2}{n_a}}$$

$$\sigma_2 = \pm \sqrt{\frac{\sum(v - \bar{v}_j)^2}{n_b}} \quad \text{---}$$

$$\sigma_c \text{ ---}$$

Mas em vez destes erros usa-se em geral uma estimativa geral, que é o erro dentro, calculado pela fórmula: (veja fórmula 7-1):

$$\begin{aligned} \sigma_b &= \pm \sqrt{\frac{\sum(u-\bar{u})^2}{m \cdot n_p}} = \pm \sqrt{\frac{\sum(u-\bar{u}_a)^2 + \sum(u-\bar{u}_b)^2}{m \cdot n_p}} \\ &= \pm \sqrt{\frac{n_p \sigma_a^2 + n_p \sigma_b^2 \dots}{m \cdot n_p}} = \pm \sqrt{\frac{\sigma_a^2 + \sigma_b^2 \dots}{m}} \end{aligned}$$

Analisando as diferenças entre médias parciais e a média geral, temos:

$$D = \frac{\bar{u} - \bar{u}}{\sigma \text{ dif.}} = \frac{\bar{u} - \bar{u}}{\sigma : \sqrt{N_p}} = \frac{\bar{u} - \bar{u}}{\sigma} \cdot \sqrt{N_p} \dots$$

Que valor de σ deveremos usar? Poderemos empregar o erro de cada média, por exemplo:

$$\left. \begin{aligned} D &= \frac{\bar{u}_a - \bar{u}}{\sigma_a} \cdot \sqrt{N_p} \\ D &= \frac{\bar{u}_b - \bar{u}}{\sigma_b} \cdot \sqrt{N_p} \end{aligned} \right\}$$

e comparar este valor com os limites da respectiva distribuição de STUDENT com n_p graus de liberdade.

Porém, em vez disso, emprega-se geralmente o erro dentro, que por sua vez é calculado de $m \cdot n_p$ graus de liberdade:

$$D = \frac{\bar{u} - \bar{u}}{\sigma_p} \cdot \sqrt{N_p} \dots$$

Este termo segue uma distribuição de STUDENT, com $m \cdot np$ graus de liberdade.

Substituindo σ_1 por σ_D , não somente usamos um outro valor do erro, mas atribuímos ao desvio relativo um grau de liberdade muito mais elevado, isto é, $(m) \cdot np$ em vez de np . O efeito desta substituição relativo ao limite de probabilidade é muito grande, pois em geral np é bem pequeno, raramente mais do que 5, e o valor de $(m) \cdot np$ chega muitas vezes perto ou além de 30. Dêste modo e conforme explicado acima, substituímos os limites de uma distribuição de STUDENT com n_1 igual ou menor do que 5 pelos limites da distribuição de GAUSS.

Antes de substituir o valor individual σ_1 , por uma média balanceada σ_D , seria imprescindível provar que a substituição é permissível. Podemos comparar os dois erros, determinando o quociente

$$\frac{\sigma_1}{\sigma_2} \text{ com } \frac{n_1 = np}{n_2 = m \cdot np}$$

e comparando-o com os limites da distribuição de FISHER correspondente.

Este teste em geral é de pouca confiança quando n_1 é muito pequeno. Quando σ_1 fôr estatisticamente diferente de σ_D o que não é raro, a substituição não pode ser feita. Mas quando a comparação não indica diferenças significantes, podemos ainda ficar na dúvida por causa do número pequeno de graus de liberdade. Nestes casos recomendo fazer a substituição e usar em vez de σ_1 o valor σ_D que provavelmente é uma melhor estimativa do erro standard, porém continuar a usar o grau de liberdade que se refere a σ_1 , isto é np .

Não é nossa intenção entrar aqui em mais detalhes, uma vez que pretendemos tratar dos testes estatísticos em outra publicação. Os exemplos servem para mostrar que os técnicos devem tomar cuidado na aplicação dos limites de probabilidade na análise de casos complicados, para evitar que eles empreguem conceitos e limites errados.

As nossas considerações finais indicam também em que direção novos trabalhos são necessários. A preparação de uma tabela com os limites bilaterais para as distribuições de FISHER é uma tarefa laboriosa, porém não envolve mais um estudo de problemas fundamentais. O principal problema de ordem teórica é a combinação de diferentes distribuições para

os casos nos quais devemos, com todo rigor, comparar variáveis que seguem cada um tipo diferente de distribuição de acaso. O teste de BEHREND'S é um início nesta direção.

RESUMO

1) Chamamos um desvio relativo simples o quociente de um desvio, isto é, de uma diferença entre uma variável e sua média ou outro valor ideal, e o seu erro standard.

$$D = \frac{v - \bar{v}}{\sigma} \quad \text{ou} \quad D = \frac{v - v_0}{\sigma}$$

Num desvio composto nós reunimos vários desvios de acôrdo com a equação:

$$D = \pm \sqrt{\frac{\sum (v - \bar{v})^2}{n_1}} : \sigma_0 = \frac{\sigma_1}{\sigma_0}$$

Todo desvio relativo é caracterizado por dois graus de liberdade (número de variáveis livres) que indicam de quantas observações foi calculado o numerador (grau de liberdade n_1 ou simplesmente n_2) e o denominador (grau de liberdade n_1 ou simplesmente n_2).

2) Explicamos em detalhe que a chamada distribuição normal ou de GAUSS é apenas um caso especial que nós encontramos quando o erro standard do dividendo do desvio relativo é calculado de um número bem grande de observações ou determinado por uma fórmula teórica. Para provar este ponto foi demonstrado que a distribuição de GAUSS pode ser derivada da distribuição binomial quando o expoente desta torna-se igual a infinito (Fig.1).

3) Assim torna-se evidente que um estudo detalhado da variação do erro standard é necessário. Mostramos rapidamente que, depois de tentativas preliminares de LEXIS e HELMERT, a solução foi achada pelos estatísticos da escola londrina: KARL PEARSON, o autor anônimo conhecido pelo nome de STUDENT e finalmente R. A. FISHER.

4) Devemos hoje distinguir quatro tipos diferentes de dis-

tribuições de acaso dos desvios relativos, em dependência de combinação dos graus de liberdade n_1 e n_2 .

Distribuição de:

<i>Fisher</i>	$1 < n_1 < infinito$	$1 < n_2 < infinito$	(formula 3-1)
<i>Pearson</i>	$1 < n_1 < infinito$	$n_2 = infinito$	(formula 3-2)
<i>Student</i>	$n_1 = 1$	$1 < n_2 < infinito$	(formula 3-3)
<i>Gauss</i>	$n_1 = 1$	$n_2 = infinito$	(formula 3-4)

As formas das curvas (Fig. 2) e as fórmulas matemáticas dos quatro tipos de distribuição são amplamente discutidas, bem como os valores das suas constantes e de ordenadas especiais.

5) As distribuições de GAUSS e de STUDENT (Figs. 2 e 5) que correspondem a variação de desvios simples são sempre simétricas e atingem o seu máximo para a abscissa $D = 0$, sendo o valor da ordenada correspondente igual ao valor da constante da distribuição, k_1 e k_2 respectivamente.

6) As distribuições de PEARSON e FISHER (Fig. 2) correspondentes à variação de desvios compostos, são descontínuas para o valor $D = 0$, existindo sempre duas curvas isoladas, uma à direita e outra à esquerda do valor zero da abscissa.

As curvas são assimétricas (Figs. 6 a 9), tornando-se mais e mais simétricas para os valores elevados dos graus de liberdade.

7) A natureza dos limites de probabilidade é discutida. Explicamos porque usam-se em geral os limites bilaterais para as distribuições de STUDENT e GAUSS e os limites unilaterais superiores para as distribuições de PEARSON e FISHER (Figs. 3 e 4).

Para o cálculo dos limites deve-se então lembrar que o desvio simples, $D = (\sqrt{v} - \bar{v}) : \sigma$ tem o sinal positivo ou negativo, de modo que é em geral necessário determinar os limites bilaterais em ambos os lados da curva (GAUSS e STUDENT).

Os desvios relativos compostos da forma $D = \sigma_1 : \sigma_2$ não têm sinal determinado, devendo desprezar-se os sinais. Em geral consideramos apenas o caso σ_1 ser maior do que σ_2 e os limites se determinam apenas na extremidade da curva que corresponde a valores maiores do que 1. (Limites unilaterais superiores das distribuições de PEARSON e FISHER).

Quando a natureza dos dados indica a possibilidade de aparecerem tanto valores de σ_1 maiores como menores do que σ_2 , devemos usar os limites bilaterais, correspondendo os limites unilaterais de 5%, 1% e 0,1% de probabilidade, correspondendo a limites bilaterais de 10%, 2% e 0,2%.

8) As relações matemáticas das fórmulas das quatro distribuições são amplamente discutidas, como também a sua transformação de uma para outra quando fazemos as necessárias alterações nos graus de liberdade. Estas transformações provam matematicamente que todas as quatro distribuições de acaso formam um conjunto.

Foi demonstrado matematicamente que a fórmula das distribuições de FISHER representa o caso geral de variação de acaso de um desvio relativo, se nós extendermos a sua definição desde $nf_1 = 1$ até infinito e desde $nf_2 = 1$ até infinito.

9) Existe apenas uma distribuição de GAUSS; podemos calcular uma curva para cada combinação imaginável de graus de liberdade para as outras três distribuições. Porém, é matematicamente evidente que nos aproximamos a distribuições limitantes quando os valores dos graus de liberdade se aproximam ao valor infinito.

Partindo de fórmulas com área unidade e usando o erro standard como unidade da abcissa, chegamos às seguintes transformações:

a) A distribuição de STUDENT (Fig. 5) passa a distribuição de GAUSS quando o grau de liberdade n_2 se aproxima ao valor infinito. Como aproximação ao infinito, suficiente na prática, podemos aceitar valores maiores do que $n_2 = 30$.

b) A distribuição de PEARSON (Fig. 6) passa para uma de GAUSS com média zero e erro standard unidade quando n_1 é igual a 1.

Quando de outro lado, n_1 torna-se muito grande, a distribuição de PEARSON podia ser substituída por uma distribuição modificada de GAUSS, com média igual a 1 e unidade da abcissa igual a $1 : \sqrt{2 n_1}$. Para fins práticos, valores de n_1 maiores do que 30 são em geral uma aproximação suficiente ao infinito.

c) Os limites da distribuição de FISHER são um pouco mais difíceis para definir.

I) Em primeiro lugar foram estudadas as distribuições com $n_1 = n_2 = n$ e verificamos (Figs. 7 e 8) que aproximamo-nos a uma distribuição transformada de GAUSS com média 1 e erro standard $1 : \sqrt{n}$, quando o valor cresce até o infinito. Como aproximação satisfatória podemos considerar $n_1 = n_2 = 100$, ou já $n_1 = n_2 = 50$ (Fig. 8)

II) Quando n_1 e n_2 diferem (Fig. 9) podemos distinguir dois casos:

Se n_1 é pequeno e n_2 maior do que 100 podemos substituir a distribuição de FISHER pela distribuição correspondente de PEARSON. (Fig. 9, parte superior).

Se porém n_1 é maior do que 50 e n_2 maior do que 100, ou vice-versa, atingimos uma distribuição modificada de GAUSS com média 1 e erro standard

$$1 : \sqrt{\frac{2n_1 n_2}{n_1 + n_2}}$$

10) As definições matemáticas e os limites de probabilidade para as diferentes distribuições de acaso são dadas em geral na literatura em formas bem diversas, usando-se diferentes sistemas de abcissas.

Com referência às distribuições de FISHER, foi usado por este autor, inicialmente, o logaritmo natural do desvio relativo, como abcissa. SNEDECOR (1937) emprega o quadrado dos desvios relativos e BRIEGER (1937) o desvio relativo próprio.

As distribuições de PEARSON são empregadas para o X^2 teste de PEARSON e FISHER, usando como abcissa os valores de

$$X^2 = D^2 \cdot n_1$$

Foi exposto o meu ponto de vista, que estas desigualdades trazem desvantagens na aplicação dos testes, pois atribui-se um pêso diferente aos números analisados em cada teste, que são somas de desvios quadrados no X^2 teste, somas des desvios quadrados divididos pelo grau de liberdade ou varianças no F-teste de SNEDECOR, desvios simples no t-teste de STUDENT, etc..

Uma tábua dos limites de probabilidade de desvios relativos foi publicada por mim (BRIEGER 1937) e uma tábua mais extensa será publicada em breve, contendo os limites unilaterais e bilaterais, tanto para as distribuições de STUDENT como de FISHER.

11) Num capítulo final são discutidas várias complicações que podem surgir na análise. Entre elas quero apenas citar alguns problemas.

a) Quando comparamos o desvio de um valor e sua média, deveríamos corretamente empregar também os erros de ambos estes valores:

$$D = \frac{v - \bar{v}}{\sqrt{\sigma^2 + \sigma_{\bar{v}}^2}}$$

Mas não podemos aqui imediatamente aplicar os limites de qualquer das distribuições do acaso discutidas acima. Em geral a variação de v , medida por σ , segue uma distribuição de STUDENT e a variação da média \bar{v} segue uma distribuição de GAUSS. O problema a ser solucionado é, como reunir os limites destas distribuições num só teste. A solução prática do caso é de considerar a média como uma constante, e aplicar diretamente os limites de probabilidade das distribuições de STUDENT com o grau de liberdade do erro σ . Mas este é apenas uma solução prática. O problema mesmo é, em parte, solucionado pelo teste de BEHREND'S.

b) Um outro problema se apresenta no curso dos métodos chamados "analysis of variance" ou decomposição do erro.

Supomos que nós queremos comparar uma média parcial \bar{v}_a com a média geral \bar{v} . Mas podemos calcular o erro desta média parcial, por dois processos, ou partindo do erro individual σ_a ou do erro "dentro" σ_D que é, como explicado acima, uma média balanceada de todos os m erros individuais. O emprêgo deste último garante um teste mais satisfatório e severo, pois é baseado sempre num grau de liberdade bastante elevado.

Teremos que aplicar dois testes em seguida:

Em primeiro lugar devemos decidir se o erro σ_a difere do erro dentro:

$$D = \frac{\sigma_a}{\sigma_D} \quad \frac{n_1 = n_p}{n_2 = m \cdot n_p}$$

Se este teste for significativo, uma substituição de σ_a pelo σ_D não será admissível. Mas mesmo quando o resultado for insignificante, ainda não temos certeza sobre a identidade dos dois erros, pois pode ser que a diferença entre eles é pequena e os graus de liberdade não são suficientes para permitir o reconhecimento desta diferença como significativa.

Podemos então substituímos σ_a por σ_D de modo que $n_2 = m : np$:

$$D = \frac{\bar{u}_a - \bar{u}}{\sigma_a} \sqrt{N_p} \quad \frac{n_1 = 1}{n_2 = np}$$

) passa para

$$D = \frac{\bar{u}_a - \bar{u}}{\sigma_D} \sqrt{N_p} \quad \frac{n_1 = 1}{n_2 = m \cdot np}$$

Mas como podemos incluir neste último teste uma apreciação das nossas dúvidas sobre o teste anterior $\sigma_a; \sigma_D$?

A melhor solução prática me parece fazer uso da determinação de σ_D , que é provavelmente mais exata do que σ_a , mas usar os graus de liberdade do teste simples: $np = 1 / 1/2 = np$ para deixar margem para as nossas dúvidas sobre a igualdade de σ_a a σ_D .

Estes dois exemplos devem ser suficientes para demonstrar que apesar dos grandes progressos que nós podíamos registrar na teoria da variação do acaso, ainda existem problemas importantes a serem solucionados.

ABSTRACT

1) The present paper deals with the mathematical basis and the relations of the different chance distributions. It is shown that the concepts of classical statistics may only be applied correctly when dealing with illimited populations where the number of variables is so large that it may be considered as infinite. After the attempts of LEXIS and HELMERT, a partial solution was found by KARL PEARSON and by STUDENT, until finally R. A. FISHER gave the general solution, solving the problem of statistical analysis in a general form and determining the chance distribution in small samples.

2) As a basis for the formulas, I am using always the relative deviate, which may be determined in two ways:

the simple relative deviate:

$$D = \frac{v - \bar{v}}{\sigma}$$

the compound relative deviate:

$$D = \pm \sqrt{\frac{\sum(v - \bar{v})^2}{n_1}} : \sigma_0 \quad \text{ou} \quad D = \frac{\sigma_1}{\sigma_2}$$

3) The deviates are always defined by two degrees of freedom, n_1 for the dividend and n_2 for the divisor. According to the values combined in any given case, we may distinguish four basic chance distributions which we shall call according to the respective authors: the distributions of GAUSS, STUDENT, PEARSON and FISHER.

The mathematical definition and the corresponding degrees of freedom are given both in formulae 3-1 to 3-4 on pg. and in the lower part of Fig. 2. The upper part of Fig. 2. represents grafically these four distributions. The equations and the forma of the corresponding curves are discussed in detail.

4) The main differences between the simple and the compound relative deviate are discussed:

a) Simple deviates have always a definite signe and are either positive or negative, according to the signe of the numerator. Correspondingly the distributions of GAUSS and STUDENT are symetrical with regards to the abscissa zero and extend on both sides of it untill plus and minus infinite. Compound deviates on the other side, have no definite sign, since the numerator is a square root. The distributions of PEARSON and FISHER, accordingly, are discontinuous for the value zero and we obtain two identical and independent curves which go from zero to plus infinite, resp. from zero to minus infinite.

b) Secondly when studying simple deviates we admitt that both positive and negative large deviates may occur in consequence of an increase in variability. Consequently we ha-

ve to use, in the corresponding tests, bilateral limits of probability (Fig. 3).

When analysing compound deviates, we are comparing one standard error with another, which may either be an ideal value or at least a better estimate. Admitting that only an increase of variability may occur, we apply in tests, based on PEARSON's or FISHER's distributions, only the upper (superior) unilateral limit of probability (Fig. 4).

The tables thus far published, for these distributions contain the unilaeral limits only. A more complete table, including bilateral limits, has been computed by myself and is already in press.

5) Discussing the relations of the four distributions, it is shown that mathematically their formulas can be easily transformed from one to the other by changing the respective values of degrees of freedom. The application of a few principles of mathematics is sufficient, besides remembering that the distributions of PEARSON and FISHER correspond only to half a distribution of STUDENT and GAUSS.

Thus it is shown:

a) that for n_1 bigger than 30, the distribution of STUDENT is so near to that of GAUSS (or normal), to permit its substitution.

b) that for n_1 bigger than 30, the distribution of PEARSON becomes almost symetrical and may be substituted by a modified distribution of GAUSS (or normal) with mean equal to one and error standard

$$1: \sqrt{2n_1}$$

c) That the distribution of FISHER with $n_1 = n_2$ becomes more or less symetrical when both reach the limit of 50 or better still 100, and than may be substituted by a modified distribution of GAUSS with mean one and error standard.

$$1: \sqrt{n}$$

d) That the distributions of FISHER, when n_1 differs from n_2 , may be substituted either by the correspondent distribution of PEARSON, if n_1 is small and n_2 bigger than 100, or by a modified distribution of GAUSS with mean unity and error standard equal to

$$1. \sqrt{\frac{2n_1n_2}{n_1+n_2}}$$

when n_1 goes beyond 50 and n_2 beyond 100 ou vice versa.

6) The formulas, generally given in the litterature to characterize the different distributions are far from being uniform and use differents measures for the abscissa.

Thus in the tests for FISHER's distribution, the natural logarithm for the deviate were used initially (FISER's z-test). Later on SNEDECOR (1937) recommended the square of the deviate (F-test) and BRIEGER (1937) the deviate itself (ψ -test).

In the X^2 test, based on PEARSON's distribution, one generally uses the square of the compound deviate, multiplied by the degree of freedom n_1 .

The t-test, based on STUDENT's distribution, finally makes use of the simple deviate itself.

The inevitabal algebraic consequences of this variation of units of emasure is, that the severity and thus the statistical efficiency of the tests is not comparable. Decimals in the t-test and ψ -test correspond to hundreds in the F-test and to almost anything, depending upon the values of n_1 , in the X^2 -test.

7) In the last chapter a few rather complicated problems are discussed, which can be solved with approximation in practical tests, but wich are still unsolved from the theoretical point of view. We shall mention here only two of the questions raised:

a) Analyzing the difference between a variable and its mean (or of a partial mean and a general mean), only the standard error of the first term is used generally, considering the other as a constant:

$$D = \frac{v - \bar{v}}{\sigma_a}$$

However with more justification both terms may be considered as variable und thus one should apply the formula :

$$D = \frac{v - \bar{v}}{\sqrt{\sigma^2 - \sigma_b^2}}$$

The first mentioned simple value of D should follow a distribution of STUDENT and its analysis thus does not present any difficulties. But in the second term we combine the term, v , with standard error σ which should follow STUDENT's distribution and the mean, \bar{v} , with standard error σ_v which generally will follow the distribution of GAUSS. How shall we combine the requirements of those two distributions simultaneously? BEHREND's test seems to give a solution, which however is not very easy to apply and which is not sufficient when the second term follows also a distribution of STUDENT, but with different degree of freedom.

b) The second problem arises in connection with the analysis of variance in its most simple form, i. e. the test "within-between". If we want to compare by a t -test the partial mean of one sample \bar{v}_a with the general mean \bar{v} of the whole experiment, we must decide whether we should use standard error of this sample σ_a , based on np degrees of freedom or the error "within" σ_D which is a balanced mean value of all the m individual sample errors. At the same time we have an alternative choice with regards to the degree of freedom:

$$D = \frac{\bar{v}_a - \bar{v}}{\sigma_a} \sqrt{N_p} \quad \frac{n_1 = 1}{n_2 = np} \quad \text{ou}$$

$$D = \frac{\bar{v}_a - \bar{v}}{\sigma_D} \sqrt{N_p} \quad \frac{n_1 = 1}{n_2 = m \cdot np}$$

Thus it is evident that the use of the value σ_D not only alters the value of the relative deviate D , but also the limits of probability to be applied which depend upon the degree of freedom. However we must justify the substitution of the partial error σ_a by the error "within" σ_D and this should be done by determining whether the value σ_a : σ_D is due to chance only, i. e. that there is really no difference between the two errors from a statistical point of view. The necessary test however:

$$D = \frac{\sigma_a}{\sigma_D} \quad \frac{n_1 = np}{n_2 = m \cdot np}$$

generally does not allow a very decisive answer since the degree of freedom np is in most cases small.

Whenever there is some reason to doubt whether the substitution is really justified, it seems to me reasonable to use the probably better estimate σ_D , instead of the individual sample error σ_a , while at the same time make allowances for doubts by not substituting the degrees of freedom:

$$D = \frac{\bar{u}_a - \bar{u}}{\sigma_D} \sqrt{N_p} \quad \frac{n_1 = 1}{n_2 = n_p}$$

A more complete formula naturally would be the following:

$$D = \frac{\bar{u}_a - \bar{u}}{\sigma_a \left(\frac{\sigma_a}{\sigma_D} \right)} \sqrt{N_p} \quad \frac{n_1 = 1}{n_2 = \dots}$$

c) These two examples should be sufficient to show that there are still important theoretical problems to be solved, in spite of the really very considerable progress achieved with regards to theory and methods of analysis of simple and compound relative deviates from uniform small or large, but always limited samples.

REFERÊNCIAS

- 1 — BRIEGER, F. G. — 1937 — Tábuas e Fórmulas para Estatística, 46 pgs. — Cia. Melhoramentos de S. Paulo — S. Paulo.
- 2 — FISHER, R. A. — 1941 — Statistical Methods for Research Workers — 8th Ed. Oliver and Boyd — London — 334 pgs.
- 3 — FISHER, R. A. and FRANK YATES — 1943 — Statistical Tables, for Biological, Agricultural and Medical Research — Second Edition — Oliver and Boyd Ltd. — London — 98 pgs.
- 4 — KENNEY, J. F. — 1939 — Mathematics of Statistics — Van Nostrand Company — New York — Vol. 1, 248 pgs — Vol. 2, 202 pgs.
- 5 — SNEDECOR G. W. — 1937 — Statistical Methods — Collegiate Press — Ames, Iowa.
- 6 — YULE, G. U. and M. G. KENDALL — 1940 — An Introduction to the Theory of Statistics — Charles Griffin 12th — 570 pgs.

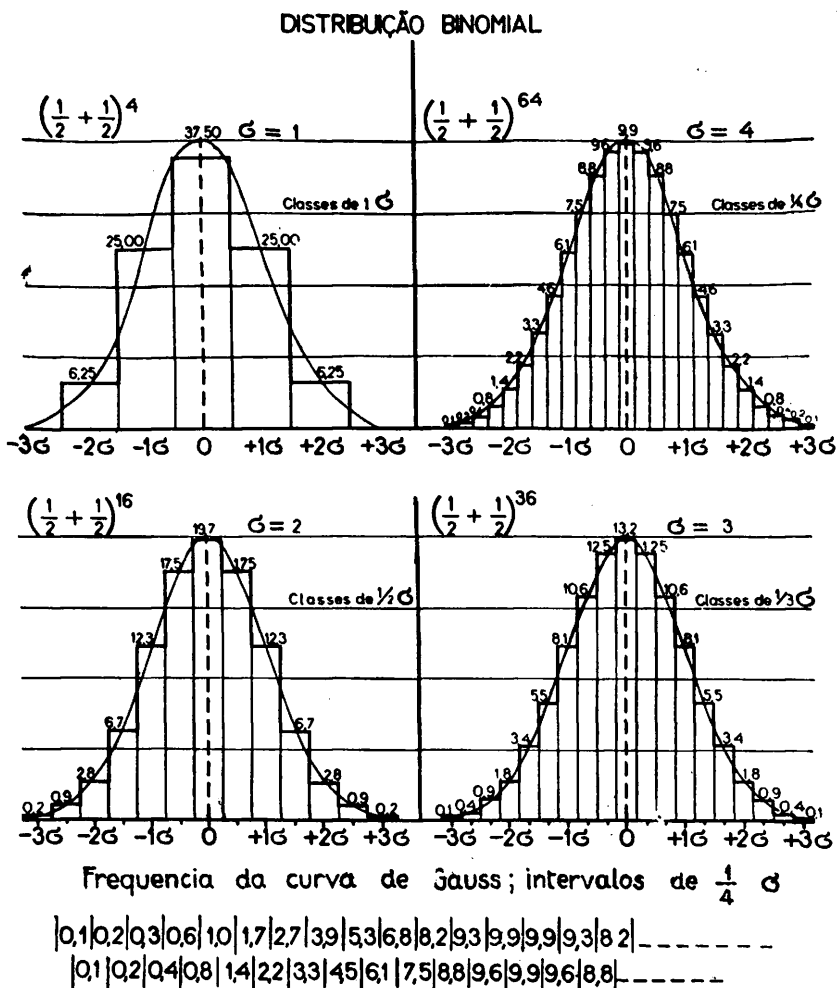
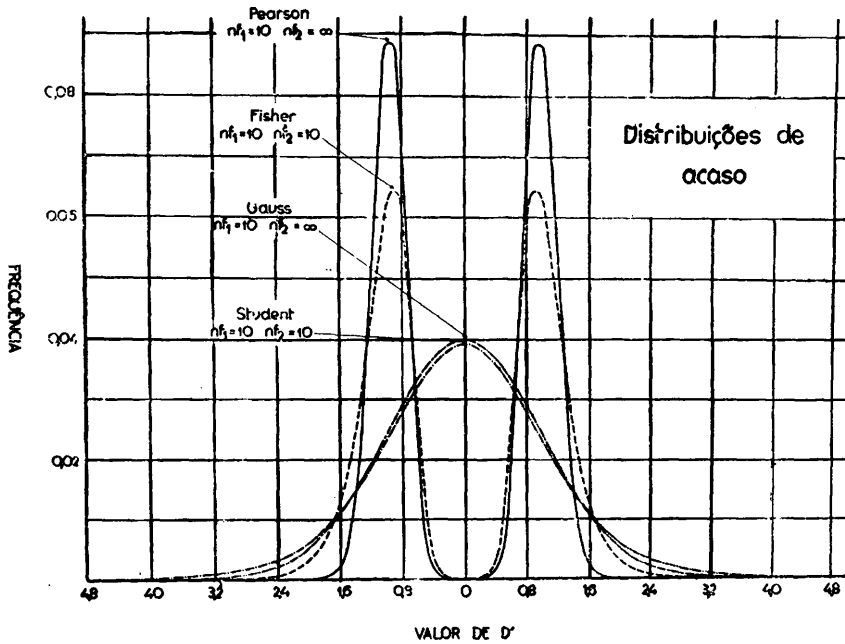


Fig. 1



Fisher:

$$y = k_1 \cdot D^{n_1-1} \cdot \left(1 + \frac{n_1}{n_2} D^2\right)^{-\frac{n_1+n_2}{2}}$$

Pearson:

$$y = k_2 \cdot D^{n_1-1} \cdot e^{-\frac{n_1}{2} D^2}$$

Student:

$$y = k_3 \cdot \left(1 + \frac{1}{n_2} D^2\right)^{-\frac{n_2+1}{2}}$$

Gauss:

$$y = k_4 \cdot e^{-\frac{1}{2} D^2}$$

Grãos de liberdade

$$\left. \begin{array}{l} 1 < n_1 < \infty \\ 1 < n_2 < \infty \end{array} \right\} D = \frac{\sigma_1}{\sigma_2}$$

$$\left. \begin{array}{l} 1 < n_1 < \infty \\ n_2 = \infty \end{array} \right\}$$

$$\left. \begin{array}{l} n_1 = 1 \\ 1 < n_2 < \infty \end{array} \right\}$$

$$\left. \begin{array}{l} n_1 = 1 \\ n_2 = \infty \end{array} \right\}$$

$$D = \frac{\sigma_1}{\sigma}$$

Fig. 2

Distribuição de Gauss

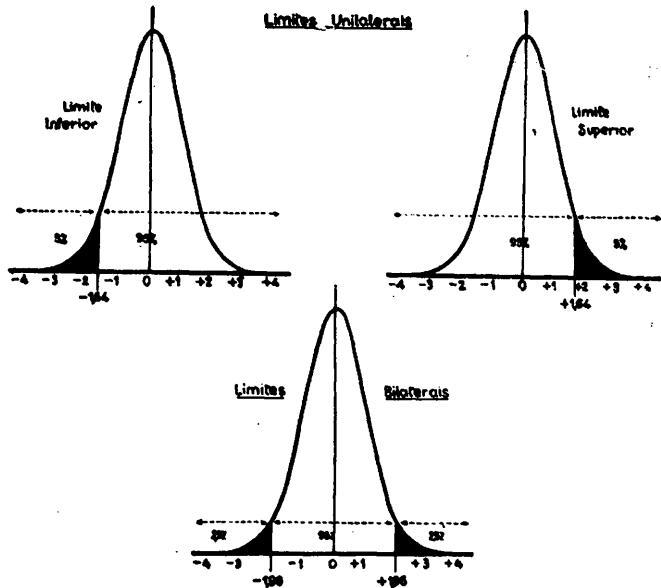


Fig. 3

Distribuição de Fisher: $nf_1 = nf_2 = 40$

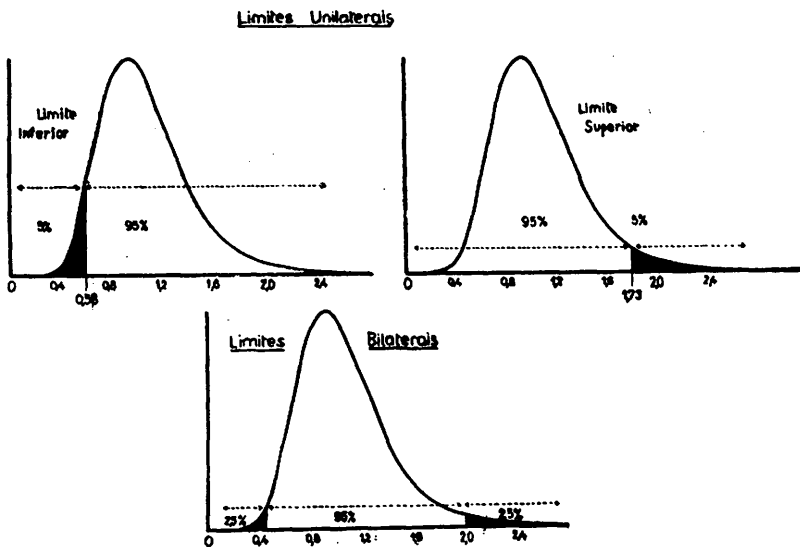
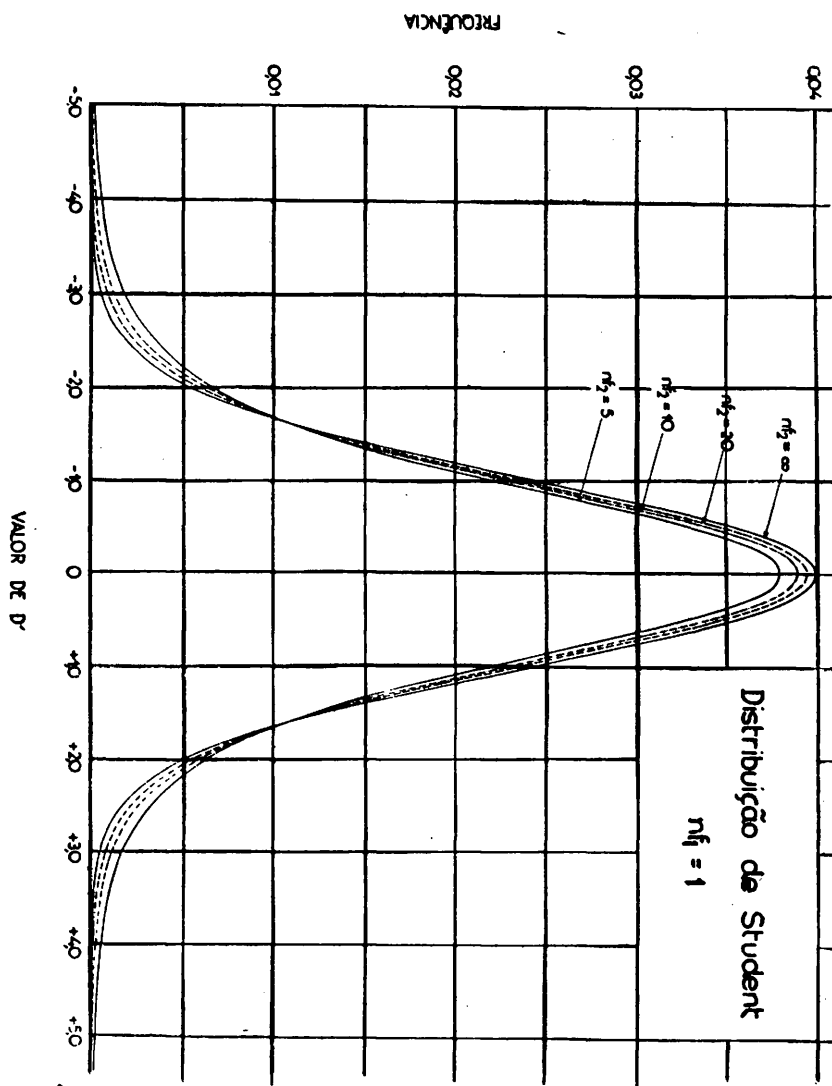


Fig. 4



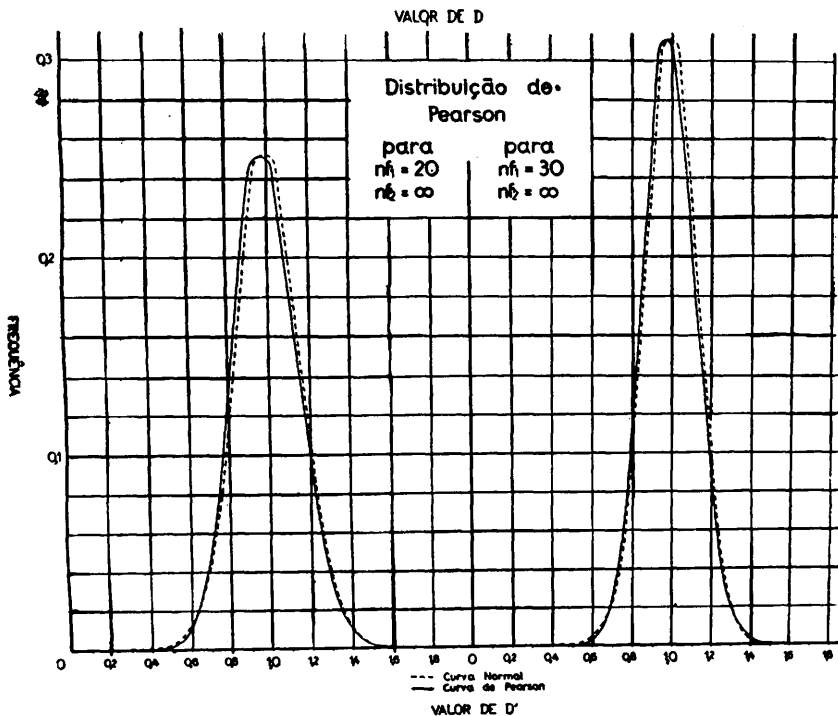
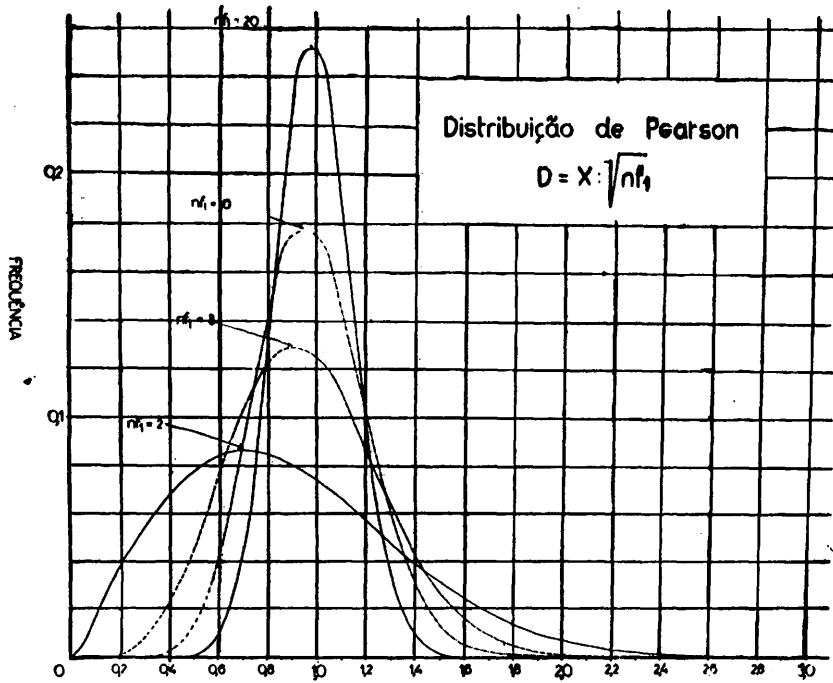


Fig. 6

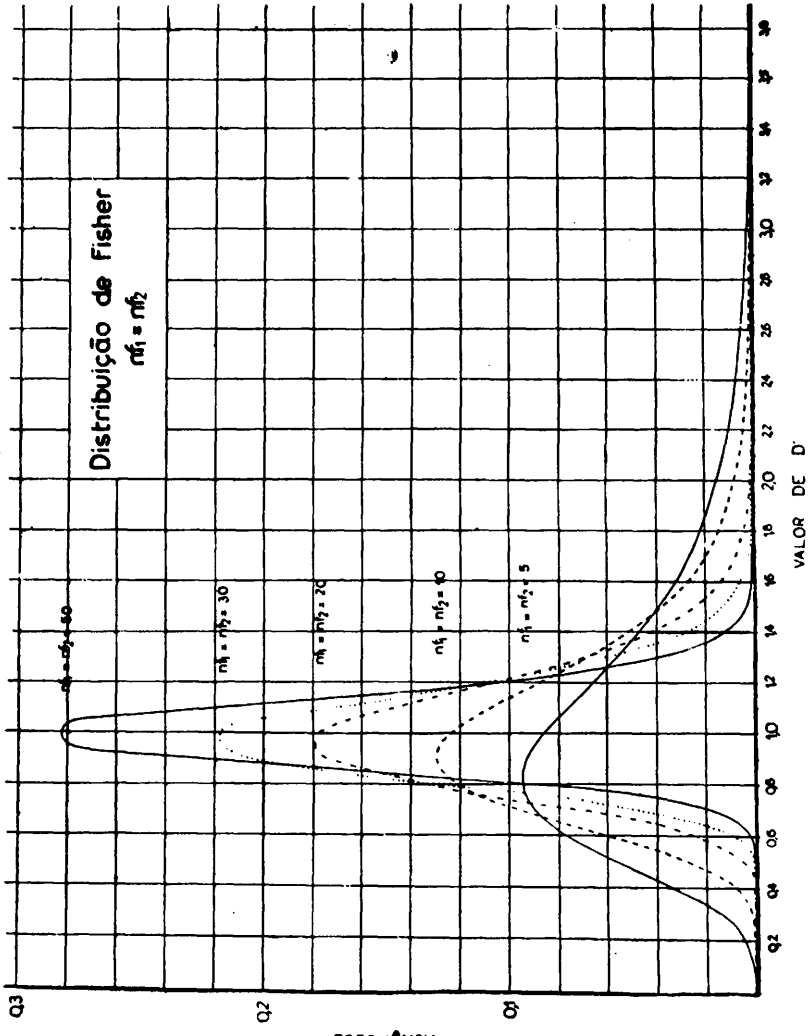


Fig. 7

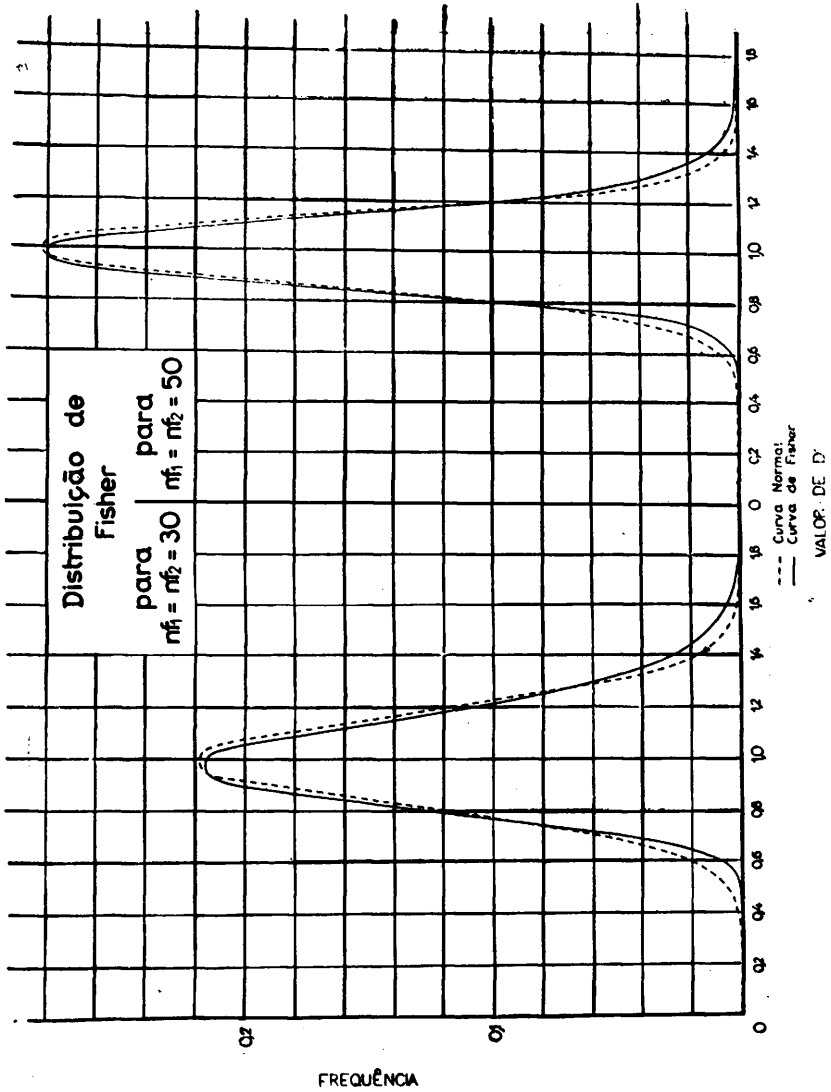


Fig. 8

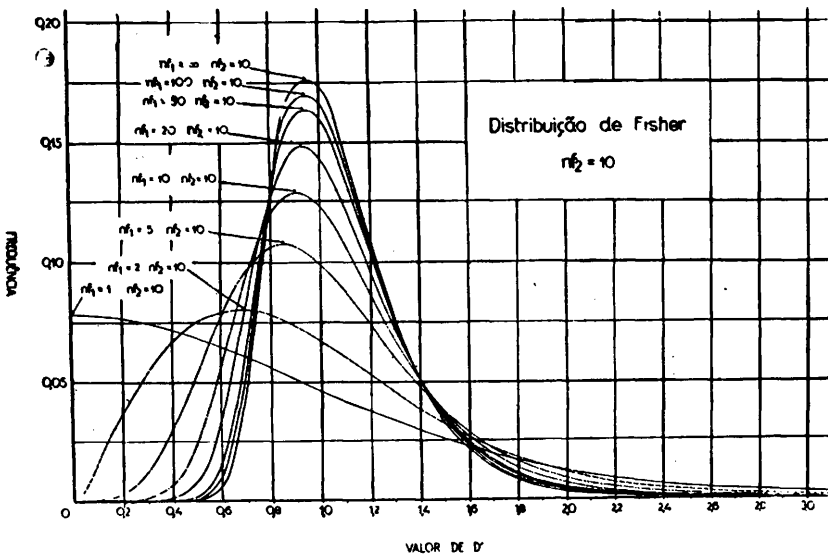
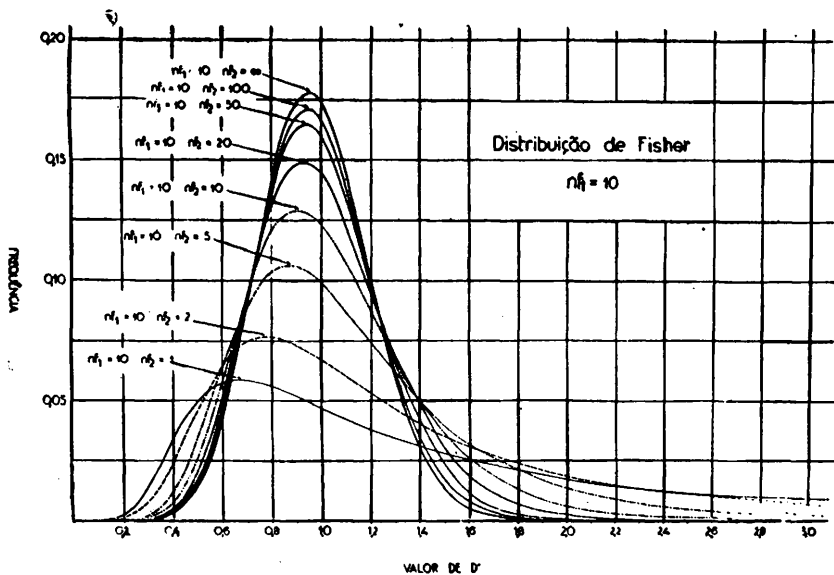


Fig. 9

NOTA

O manuscrito do presente trabalho serviu como base a várias conferências, realizadas no Ministério de Agricultura no Rio de Janeiro e na Escola Superior Agrícola de Viçosa, em Setembro de 1944. Por razões técnicas atrasou-se a impressão.

CORREÇÕES

Infelizmente escaparam ao autor alguns enganos nas fórmulas :

Pg. 343 — Em baixo do termo: $\lim y$ [Fisher] deve constar: $n_2 = \text{inf.}$ em vez de $n_1 = \text{inf.}$

Pg. 348 — Em baixo do termo: $\lim y$ [Fisher] deve constar: $n_2 = \text{inf.}$ em vez de $n_1 = 1$.

Pg. 349 — Na fórmula do termo k_2 , no início da página, deve-se ler no dividendo n_1 , elevado $n_1 : 2$, em vez de n_2 .

Pg. 360 — Deve-se ler na segunda linha das fórmulas: "Substituímos agora D por $X : n, \frac{1}{2}$, em vez de $X n, \frac{1}{2}$."

Pg. 365 — No último termo da fórmula 7-8 falta o fator σ antes da raiz quadrada no divisor.

Pg. 371 — Na fórmula, no centro da página, deve-se ler $\sigma_1 : \sigma D$, em vez de $\sigma_1 : \sigma_2$.

Pg. 380 — Na primeira fórmula da página deve-se substituir o sinal de multiplicação pelo de divisão.

Na última fórmula deve-se substituir o sinal de multiplicação pelo de igualdade.