

## Por que a fórmula para a variância amostral tem $n - 1$ no seu denominador

Suponhamos que alguém colocou bolas idênticas numa urna, sendo que cada bola carrega um dos números  $k_1, k_2, \dots, k_M$ . Denotaremos por  $p_i$  a proporção das bolas da urna que carregam o número  $k_i$ . Com isso, a seguinte tabela

$k_1$	$k_2$	$\dots\dots\dots$	$k_{M-1}$	$k_M$
$p_1$	$p_2$	$\dots\dots\dots$	$p_{M-1}$	$p_M$

é exatamente aquilo ao que chamamos por “distribuição de frequência relativa por atributo “numero carregado” da população de bolas na urna”. Esse distribuição será chamada no que se segue por “distribuição populacional”.

Recordo, para as necessidades dos argumentos a vir, que a quantia

$$k_1p_1 + k_2p_2 + \dots + k_{M-1}p_{M-1} + k_Mp_M \quad (1)$$

calculada com valores da distribuição populacional, chama-se *média populacional*; denotaremos essa por  $\mu$ .

Outrossim recordo, também para as necessidades futuras, que a quantia

$$(k_1 - \mu)^2p_1 + (k_2 - \mu)^2p_2 + \dots + (k_{M-1} - \mu)^2p_{M-1} + (k_M - \mu)^2p_M \quad (2)$$

chama-se *variância populacional*; denotaremos essa por  $\sigma^2$ .

Imagine agora que não sabemos nada sobre a distribuição populacional: não sabemos nem  $M$ , nenhum dos  $k_i$  e nenhum dos  $p_i$ .

Imagine que desejamos estimar o valor de  $\sigma^2$ , quer dizer, desejamos estimar a variância populacional, e para o fim dessa estimação, somos autorizados fazer amostra simples com reposição. Isso quer dizer, que podemos retirar ao acaso uma bola da urna e observar o número carregada pela bola retirada, e podemos repetir esse procedimento tantas vezes quantas desejamos, desde que após cada retirada, a bola seja devolvida à urna, e todas as bolas sejam bem misturadas.

Então, ao fixar o número de retiradas a serem feitas e ao dentro esse por  $n$  (esse  $n$  é o que chama-se *o tamanho da amostra* a ser feita), sugerimos o seguinte procedimento: após fazer  $n$  retiradas e ver os números

$$x_1, x_2, \dots, x_n \quad (3)$$

fazer a média delas, isso é,

$$\frac{1}{n} \{x_1 + x_2 + \dots + x_n\} \quad (4)$$

e usar esse para calcular o valor da expressão

$$\frac{1}{n - 1} \{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2\} \quad (5)$$



**Fato 2:**

$$\mathbb{E}[X_i] = \mu, \quad \text{Var}[X_i] = \sigma^2, \text{ para cada } X_i \quad (10)$$

**Fato 3:** As variáveis aleatórias  $X_1, X_2, \dots, X_n$  são independentes em conjunto.

Confere que você entende bem os Fatos 1-3: O primeiro deles segue-se da definição de experimento aleatório, da construção de modelo probabilístico e da definição de variável aleatória como a expressão do resultado a vir num experimento aleatório. É importante entender que a suposição, segundo a qual não conhecemos nada sobre a distribuição populacional, não impede da mesma ser a distribuição de cada variável aleatória que definimos; é só não conhecemos a distribuição. O comentário do Fato 2 vem no mesmo sentido: apesar de não conhecer os valores da média populacional e da variância populacional, sabemos que esses valores coincidem com, respectivamente, a esperança e a variância de cada uma das variáveis aleatórias que definimos. O Fato 3 segue-se de um comentário que fiz quando expliquei o conceito de independência entre variáveis aleatórias. É suficiente que você acredite na independência. Creio que isso não é difícil a ser concebido pois as bolas retiradas são devolvidas na urna, e, portanto, o resultado de  $i$ -ésima retirada não pode depender dos resultados das retiradas anteriores.

Vamos agora introduzir notações que ajudarão a simplificar a escrita da demonstração que pretendemos fazer. Em primeiro lugar, observe que

$$X_i - \bar{X} = (X_i - \mu) - (\bar{X} - \mu) = (X_i - \mu) - \frac{(X_1 - \mu) + (X_2 - \mu) + \dots + (X_n - \mu)}{n}$$

Por isso, ao introduzir

$$Y_i = X_i - \mu, \text{ para cada } i$$

podemos reescrever nosso objetivo da seguinte maneira:

$$\text{provar que } \mathbb{E} \left[ \frac{1}{n-1} \{ (Y_1 - \bar{Y})^2 + \dots + (Y_n - \bar{Y})^2 \} \right] = \sigma^2 \quad (11)$$

onde

$$Y_1, Y_2, \dots, Y_n \text{ são variáveis aleatórias independentes em conjunto,} \quad (12)$$

$$\mathbb{E}[Y_i] = 0 \text{ e } \text{Var}[Y_i] = \sigma^2 \text{ para cada } i, \quad (13)$$

$$\bar{Y} \text{ é a notação para } \frac{1}{n} \{ Y_1 + Y_2 + \dots + Y_n \} \quad (14)$$

sendo que (12) é a consequência do Fato 3, e (13) é a consequência do Fato 2; óbvio, que poderíamos desenhar a distribuição de cada  $Y_i$  a partir da distribuição das  $X$ 's, mas não fizeram isso pois estas não serão importantes para as contas a seguir.

O segundo passo na nossa presente derivação de fatos e notações auxiliares é o fato de que

$$\mathbb{E}[Y_i^2] = \sigma^2 \quad (15)$$

Esse segue-se da expansão  $\text{Var}[Y_i] = \mathbb{E}[Y_i^2] - (\mathbb{E}[Y_i])^2$  junto com as duas relações da Eq. (13). Já no terceiro passo deduzimos que

$$\mathbb{E}[Y_i Y_j] = \mathbb{E}[Y_i] \mathbb{E}[Y_j] = 0 \times 0 = 0, \text{ para todos } i \text{ e } j \text{ diferentes entre si} \quad (16)$$

Observe que a primeira igualdade da sequencia (16) sustenta-se pela independência entre  $Y_i$  e  $Y_j$ , a qual, por sua vez, segue-se da independência das  $Y$ 's em conjunto. Essa conta é o único lugar onde usamos a independência, mas o resultado da conta é a chave principal para tudo, de modo que a ausência da independência quebraria a conta e, em sequencia, toda a demonstração.

Agora, juntamos (15) e (16) para derivar que

$$\begin{aligned} \mathbb{E} [(Y_1 + Y_2 + \dots + Y_n)^2] &= \mathbb{E} [Y_1^2 + Y_2^2 + \dots + Y_n^2 + 2Y_1 Y_2 + \dots + 2Y_{n-1} Y_n] = \\ &= (\mathbb{E}[Y_1^2] + \dots + \mathbb{E}[Y_n^2]) + (2\mathbb{E}[Y_1 Y_2] + \dots + \mathbb{E}[Y_{n-1} Y_n]) = \\ &= (\sigma^2 + \dots + \sigma^2) + (0 + \dots + 0) = n\sigma^2 \end{aligned}$$

Esse resultado será usado para fechar a demonstração em conjunto com o resultado da seguinte conta, puramente algébrica:

$$\begin{aligned} \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \\ &= (Y_1^2 - 2Y_1 \times \bar{Y} + (\bar{Y})^2) + \dots + (Y_n^2 - 2Y_n \times \bar{Y} + (\bar{Y})^2) = \\ &= (Y_1^2 + \dots + Y_n^2) - 2n\bar{Y} \times \bar{Y} + n(\bar{Y})^2 = \\ &= (Y_1^2 + \dots + Y_n^2) - n(\bar{Y})^2 = (Y_1^2 + \dots + Y_n^2) - \frac{1}{n}(Y_1 + \dots + Y_n)^2 \end{aligned}$$

Portanto,

$$\begin{aligned} \mathbb{E} \left[ \sum_{i=1}^n (Y_i - \bar{Y})^2 \right] &= \mathbb{E} [Y_1^2 + \dots + Y_n^2] - \frac{1}{n} \mathbb{E} [(Y_1 + \dots + Y_n)^2] = \\ &= \mathbb{E} [Y_1^2] + \dots + \mathbb{E} [Y_n^2] - \frac{1}{n} (n\sigma^2) = \\ &= n\sigma^2 - \sigma^2 = (n-1)\sigma^2 \end{aligned}$$

Agora ficou claro que (já passando de volta para as variáveis aleatórias  $X$ 's)

$$\mathbb{E} \left[ \frac{1}{n-1} \{ (X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \} \right] = \sigma^2 \quad (17)$$

enquanto que, se tomássemos  $n$  e vez de  $n-1$  no denominador, teríamos

$$\mathbb{E} \left[ \frac{1}{n} \{ (X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \} \right] = \frac{n-1}{n} \sigma^2 \quad (18)$$