

Every human is a carrier of one of the three genotypes AA, Aa, or aa. The genotypes are occurring with the probabilities $(1 - p)^2$, $2 * p * (1 - p)$ and p^2 whereas $0 < p < 1$ and testing of n persons yielded

- x persons had the genotype AA
- y persons had the genotype Aa
- z persons had the genotype aa

Describe the corresponding statistical model and determine the Maximum Likelihood Estimator for p .

Solution:

The likelihood function is given by [1]:

$$P(x, y, z / p) = \binom{x+y+z}{x} * (1 - p)^{2 * x} * \binom{y+z}{y} * (2 * p * (1 - p))^y * \binom{z}{z} * p^{2 * z} \dots\dots\dots (1)$$

Taking log likelihood of (1) we get,

$$\ln(P(x, y, z / p)) = \ln\left(\binom{x+y+z}{x} * (1 - p)^{2 * x} * \binom{y+z}{y} * (2 * p * (1 - p))^y * \binom{z}{z} * p^{2 * z}\right)$$

$$\ln(P(x, y, z / p)) = \ln\left(\binom{x+y+z}{x}\right) + \ln((1 - p)^{2 * x}) + \ln\left(\binom{y+z}{y}\right) + \ln((2 * p * (1 - p))^y) + \ln\left(\binom{z}{z}\right) + \ln(p^{2 * z})$$

$$\ln(P(x, y, z / p)) = \text{constant}_1 + 2 * x * \ln(1 - p) + \text{constant}_2 + y * \ln(p) + y * \ln(1 - p) + \text{constant}_3 + 2 * z * \ln(p) \dots\dots\dots (2)$$

We set the derivative equal to zero:

$$\frac{2 * z + y}{p} - \frac{y + 2 * x}{1 - p} = 0 \dots\dots\dots (3)$$

Solving equation (3) we got the value of p .

$$\begin{aligned} & \frac{2 * z + y}{p} - \frac{y + 2 * x}{1 - p} = 0 \\ \rightarrow & \frac{(1-p)*(2*z+y)-p*(y+2*x)}{p*(1-p)} = 0 \\ \rightarrow & 2 * z - 2 * z * p + y - y * p - y * p - 2 * x * p = 0 \\ \rightarrow & 2 * z + y - 2 * z * p - 2 * y * p - 2 * p * x = 0 \\ \rightarrow & (2 * z + y) - p * (2 * z + 2 * y + 2 * x) = 0 \\ \rightarrow & p = \frac{2 * z + y}{2 * x + 2 * y + 2 * z} \end{aligned}$$

The corresponding statistical model is “multinomial distribution model”. An extension of the binomial distribution is the multinomial distribution. The multinomial distribution is used to simulate the results of n experiments, where each trial's outcome has a categorical distribution [3].

Exactly one of the fixed finite number k of possible results with probabilities p_1, p_2, \dots, p_k (here $p_i \geq 0$ for $i = 1, \dots, k$ and $\sum_{i=1}^k p_i = 1$), and there are n independent trials. Next, the random variable X_i indicates the number of times outcome number i was observed over the n experiments. Then $X = (X_1, X_2, \dots, X_k)$ follows a multinomial distribution with the parameters n and p . Where $p = (p_1, p_2, \dots, p_k)$ [2].

The PMF of the multinomial distribution is given by

$$P (X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \frac{n!}{x_1! x_2! \dots x_k!} * p_1^{x_1} * p_2^{x_2} \dots p_k^{x_k}$$

with, $\sum_{i=1}^k x_i = n$, and $\sum_{i=1}^k p_i = 1$

Reference:

- [1]. <https://math.mit.edu/~dav/05.dir/class10-prep.pdf>
- [2]. Sinharay, Sandip. "Discrete Probability Distributions." (2010): 132-134.
- [3]. Multinomial distribution, https://en.wikipedia.org/wiki/Multinomial_distribution