

Exercise 2

A student test consists of N questions, each one of which has k possible answers, among which the student is requested to select the only correct one ($N, k > 0$).

Jack tries his luck by making blind, independent guesses at each question.

- 1) Find the expression for the PMF P_n of the number n of correct answers in the whole test
- 2) Discuss possible approximations of P_n when $N = 80$, $k = 5$, and compute an approximated value for the probability that Jack answers 40 or more questions correctly.
- 3) (For generic N and k values): assume that each correct answer earns Jack +1 points, and each wrong answer earns it -1 points. What is the mean value of Jack's score at the test s ?
- 4) What is the score that one should get for a wrong answer, in order for the mean value of s to be equal to zero?

Now, assume that Jack has some knowledge in the subject matter of the test. He (thinks he) *knows* the answer to M questions ($0 \leq M \leq N$), which he sets apart in advance and answers to independently, and he *guesses* the remaining $N - M$ ones. When he *knows* the answer, he has 90% probability of actually getting the correct answer.

- 5) Compute the expression of the probability that Jack answers correctly to exactly M questions.

Solution

- 1) This is a repeated trial problem, with a probability of success equal to $p = 1/k$. The PMF is thus the well-known binomial PMF:

$$P_n = \binom{N}{n} p^n \cdot (1-p)^{N-n}$$

- 2) A binomial can be approximated with:
 - a. A Poisson variable, if N is large and p is small. In this case, $p = 0.2$, hence this approximation incurs large errors.
 - b. A Gaussian variable, if $N \cdot p \cdot (1-p) > 10$. With $N = 80$, $k = 5$, it is $N \cdot p \cdot (1-p) = 12.8$, hence the required condition holds.

In this case, it is $P_n \approx P(n-0.5 \leq X \leq n+0.5)$, with $X \sim N(N \cdot p, N \cdot p \cdot (1-p))$, i.e.

$$X \sim N(16, 12.8). \text{ Therefore, we have } P_n \approx P\left(\frac{n-0.5-16}{\sqrt{12.8}} \leq \frac{X-16}{\sqrt{12.8}} \leq \frac{n+0.5-16}{\sqrt{12.8}}\right).$$

With $n=40$, it is

$$\begin{aligned} P\{n \geq 40\} &= 1 - \Phi\left(\frac{23.5}{\sqrt{12.8}}\right) \\ &\approx 1 - \Phi(6.57) \\ &\approx 0 \end{aligned}$$

Since notoriously $\Phi(x) \approx 0$ when $x \geq 3$. This also makes sense intuitively. Note that the answer obtained without approximations is $P\{n \geq 40\} \cong 2.06 \cdot 10^{-9}$.

- 3) The mean score for a question is $1 \cdot p + (-1) \cdot (1-p) = 2p-1$, and the sum of the means is equal to the mean of the sum. Therefore, the mean score for the whole test is $E[s] = N \cdot (2p-1)$. This tells us that the mean test score is negative if the probability of success is less than 50% at each question, which makes sense intuitively.

The same result could be obtained via a considerably longer route as follows: with n correct answers out of N , the test score will be $n - (N-n)$. Thus, the mean value of the score is

$$E[s] = \sum_{n=0}^N P_n \cdot [n - (N-n)] = 2 \sum_{n=0}^N n \cdot P_n - N \cdot \sum_{n=0}^N P_n. \text{ This can be further developed as:}$$

$$\begin{aligned} E[s] &= 2 \sum_{n=0}^N n \cdot P_n - N \cdot \sum_{n=0}^N P_n \\ &= 2 \sum_{n=1}^N n \cdot \binom{N}{n} p^n \cdot (1-p)^{N-n} - N \\ &= 2 \sum_{n=1}^N N \cdot \binom{N-1}{n-1} p^n \cdot (1-p)^{N-n} - N \\ &= 2N \cdot p \cdot \sum_{j=0}^{N-1} \binom{N-1}{j} p^j \cdot (1-p)^{(N-1)-j} - N \\ &= 2N \cdot p - N \\ &= N \cdot (2p-1) \end{aligned}$$

- 4) In order for the mean test score to be null, *some* test scores must be negative, and this can only be obtained if the score for a wrong answer is negative. Assume that a wrong answer gets $-\delta$, with $\delta > 0$. The mean test score is null if and only if the mean score of each question is null, and this happens if $1 \cdot p + (-\delta) \cdot (1-p) = 0$, i.e. $\delta = p/(1-p) = 1/(k-1)$. The same result could be obtained via the longer route, by observing that the mean test score is:

$$E[s] = \sum_{n=0}^N P_n \cdot [n - (N-n) \cdot \delta] = (1+\delta) \sum_{n=0}^N n \cdot P_n - \delta N \cdot \sum_{n=0}^N P_n$$

Based on the computations of the previous point, one clearly sees that the mean value for the test score is:

$$E[s] = (1 + \delta) \cdot N \cdot p - \delta N = N \cdot [(1 + \delta) \cdot p - \delta].$$

Now, in order to be $E[s] = 0$, we need $\delta = p/(1-p) = 1/(k-1)$.

5) The probability that Jack gets n correct answers can be written as follows:

$$\begin{aligned} P_M' &= \sum_{j=0}^M P\{(M-j) \text{ correct guesses, } j \text{ correct "known" answers}\} \\ &= \sum_{j=0}^M (P_{M-j} | Q_j) \cdot Q_j \end{aligned}$$

Where Q_j is the probability that Jack gets j correct answers among the M that he thinks he knows. Now, it is correct to assume that P_{M-j} and Q_j are independent of each other, hence the formula is:

$$\begin{aligned} P_M' &= \sum_{j=0}^M P_{M-j} \cdot Q_j \\ &= \sum_{j=0}^M \left[\binom{N-M}{M-j} p^{(M-j)} \cdot (1-p)^{(N-M)-(M-j)} \right] \cdot \left[\binom{M}{j} q^j \cdot (1-q)^{M-j} \right] \\ &= p^M \cdot (1-p)^{(N-M)} \cdot \sum_{j=0}^M \binom{N-M}{M-j} \cdot \binom{M}{j} \cdot \left[\left(\frac{q}{p} \right)^j \cdot \left(\frac{1-q}{1-p} \right)^{M-j} \right] \end{aligned}$$

Where $q = 0.9$ is the probability that Jack answers correctly to a question whose answer he thinks he knows. Note that the above summation can start from $j = \max(0, 2M - N)$, since the first binomial coefficient is null if $j < 2M - N$.