**Section 2.5:** hypergeometric distribution.
**Section 2.7:** Poisson distribution.

### 10.1. The Hypergeometric r.v.

• This distribution is best described through an example.

EXAMPLE 10.1. An urn contains $N$ balls of which $D$ are green and $N - D$ are red. We take a sample of $n$ balls, **without** replacement. Let $X$ be the number of green balls in the sample. Find the pmf of X.

SOLUTION. The sample space has $\binom{N}{n}$ equally likely outcomes. There are $\binom{D}{k}$ ways to choose $k$ balls out of a total of $D$ green balls, and there are $\binom{N-D}{n-k}$ ways to choose $(n-k)$ balls out of a total of $N - D$ red balls. So,

$$p_X(k) = \frac{\binom{D}{k}\binom{N-D}{n-k}}{\binom{N}{n}}; \quad k = \max\left\{0, n - (N-D)\right\}, \ldots, \min\left\{n, D\right\}.$$

The support of $X$ (that is, the possible values for $k$) is justified as follows:

• minimal value is the max between 0 and $\underset{\substack{\downarrow \\ \text{sample size}}}{n} - (N - \underset{\substack{\downarrow \\ \text{\# of red balls}}}{D})$

• maximal value is the min between $D$ (number of green balls) and $n$ (sample size).

DEFINITION 10.2. A r.v. $X$ has a hypergeometric distribution with parameter $N$, $D$, and $n$ if its pmf is

$$p_X(k) = \frac{\binom{D}{k}\binom{N-D}{n-k}}{\binom{N}{n}}$$

for $\max\left\{0, n - (N-D)\right\} \le k \le \min\left\{n, D\right\}$. Notation: $X \sim \mathrm{HG}(N, D, n)$.

REMARK 10.3. If $X \sim \mathrm{HG}(N, D, n)$, then $\mathbb{E}[X] = \frac{nD}{N}$ and $\mathrm{Var}(X) = \frac{nD}{N} \cdot \left(1 - \frac{D}{N}\right) \cdot \frac{N-n}{N-1}$. This can be shown by direct calculation using the definition of mean and variance, or by using indicator random variables (try!).

### 10.2. Poisson distribution

• The pmf $p(k)$ of $\mathrm{Bin}(n, p)$ is sometimes difficult to compute, especially for large $n$. We will approximate the pmf of $\mathrm{Bin}(n, p)$ under the following assumptions:
  − $n$ is large $(n \to \infty)$
  − $p$ is small, or successes are rare $(p \to 0)$
  − $np$ is of moderate size ( $p = \frac{\lambda}{n}$ for a constant $\lambda$, or $np \to \lambda$).

- Under these assumptions, for a fixed $k$ we have:

$$p(k) = \frac{n!}{k!\,(n-k)!} \cdot \left(\frac{\lambda}{n}\right)^k \cdot \left(1 - \frac{\lambda}{k}\right)^{n-k} =$$

$$= \frac{n(n-1)(n-2)\cdots(n-k+1)}{k!} \cdot \frac{\lambda^k}{n^k} \cdot \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-k}.$$

Note that :
- $\frac{n(n-1)(n-2)\cdots(n-k+1)}{n^k} = \frac{n}{n} \times \frac{n-1}{n} \times \frac{n-2}{n} \times \cdots \times \frac{n-k+1}{n} \to 1$ $(n \to \infty$, and $k$ fixed$)$
- $\left(1 - \frac{\lambda}{n}\right)^n \to e^{-\lambda}$. Recall that $\lim_{x \to \infty} \left(1 + \frac{1}{x}\right)^x = e$.
- $\left(1 - \frac{\lambda}{n}\right)^{-k} \to 1$ $(n \to \infty$, and $k$ fixed$)$

- Combining the above, we get:

$$\lim_{n \to \infty} p(k) = \lim_{n \to \infty} \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} = \frac{\lambda^k}{k!} e^{-\lambda}.$$

- This function defines the Poisson distribution:

DEFINITION. (Poisson distribution). A r.v. $X$ has the <u>Poisson distribution</u> with parameter $\lambda > 0$ if $X$ has pmf:

$$p(k) = \frac{\lambda^k}{k!} e^{-\lambda}, k = 0, 1, 2, 3, \ldots$$

Notation: $X \sim \text{Pois}(\lambda)$.

REMARK 10.4. The Poisson distribution was introduced by Poisson in 1837, along with applications in criminal lawsuits (jury decisions). However, it had not attracted much attention until a book by Bortkiewicz was published in 1898, which included a strange example: deaths by horse kicks in the Prussian army. Using the Poisson distribution, such rare events were shown to be quite regularized and predictable, which helped to popularize this distribution.

EXAMPLE 10.5. In a class of 40 students, on average, 2 students have been sick in the past. We showed in Example 9.6 of Lecture 9 that $X$, the number of sick students, is $\text{Bin}(40, 1/20)$. We then calculated $\mathbb{P}[X = 4] = 0.09012$.
Let us use Poisson approximation. Here $X \sim \text{Bin}(40, 1/20) \approx \text{Pois}(2)$ (large $n$ and small $p$). Hence, $\mathbb{P}(X = 4) \approx \frac{2^4}{4!} e^{-2} \approx 0.09022$. This is a good approximation!

REMARK 10.6. One may use Poisson approximation to the binomial distribution, that is $\text{Bin}(n, p) \approx \text{Pois}(np)$, whenever the number of independent trials is large, $p$ is small, and $\lambda = np$ is of moderate size.

As a consequence of the fact that a Poisson r.v. is the limit of Binomial r.v.s, we get:

EXAMPLE 10.7. For $X \sim \text{Pois}(\lambda)$, show that $\mathbb{E}[X] = \lambda$ and $\text{Var}(X) = \lambda$.

SOLUTION. **_Method 1:_** The pmf of $X$ is the limit of the pmfs of Binomial r.v. $X\left(n, \frac{\lambda}{n}\right)$, hence one can expect to have $\mathbb{E}[X] \approx np = \lambda$ and $\text{Var}(X) \approx np(1-p) = \lambda(1-p) \to \lambda$ (since $p \to 0$).
**_Method 2:_** (Optional reading, not covered during the lecture) Alternatively, we can directly calculate the mean and variance. First, recall that the Taylor series of an exponential function is $e^\lambda = \sum_{k=0}^\infty \frac{\lambda^k}{k!}$. Therefore,

$$\mathbb{E}[X] = \sum_{k=0}^\infty k \frac{\lambda^k}{k!} e^{-\lambda} = \lambda e^{-\lambda} \sum_{k=1}^\infty \frac{\lambda^{k-1}}{(k-1)!} = \lambda.$$

Similarly,

$$\mathbb{E}[X^2] = \sum_{k=0}^{\infty} k^2 \frac{\lambda^k}{k!} e^{-\lambda} = \lambda e^{-\lambda} \sum_{k=1}^{\infty} \frac{k\lambda^{k-1}}{(k-1)!}.$$

$$= \lambda e^{-\lambda} \sum_{k=1}^{\infty} \frac{(k-1+1)\lambda^{k-1}}{(k-1)!}$$

$$= \lambda e^{-\lambda} \left( \sum_{k=2}^{\infty} \frac{\lambda^{k-1}}{(k-2)!} + \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} \right)$$

$$= \lambda e^{-\lambda} \left( \lambda e^{\lambda} + e^{\lambda} \right) = \lambda(1 + \lambda).$$

Finally, we obtain that $\mathrm{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \lambda + \lambda^2 - \lambda^2 = \lambda$.

- Besides approximating the binomial distribution, the Poisson distribution is an appropriate model in various situations, and has numerous applications. For example,

  (1) The number of accidents on a section of a highway on a given day

  (2) The number of floods at a river over a century.

  (3) The number of $\alpha$ particles emitted from a radioactive source during a fixed period of time.

**Poisson Process.** Often, we consider the number $N(t)$ of occurrences of an event in an interval of length $t$, say $[0, t]$, where the average number of occurrences is proportional to the length of the interval, that is $\mathbb{E}[N(t)] = \lambda t$. Furthermore, the number of events happening over non-overlapping intervals are independent. It is possible to conclude from these two conditions (and some additional technical conditions) that, for each $t$, $N(t) \sim \mathrm{Pois}(\lambda t)$. The family of r.v.s $\{N(t)\}$ for different $t$ is called a Poisson process with rate $\lambda$.