



数据科学基础

Foundations of Data Science

3.3 离散概率分布

陈振宇

南京大学智能软件工程实验室

www.iselab.cn



伯努利分布

伯努利分布 (Bernoulli distribution)，又名两点分布或0-1分布，是一个离散型概率分布，为纪念瑞士科学家雅各布·伯努利而命名。伯努利试验是只有两种可能结果的单次随机试验。

【伯努利分布】 对于一个随机试验，如果它的样本空间只包含两个元素，即 $\Omega = \{0, 1\}$ ，我们总能在 Ω 上定义一个服从伯努利分布(又称0-1分布)的随机变量。它的概率分布为：

$$P(X = 1) = p, P(X = 0) = q, p + q = 1$$

对新生婴儿的性别进行登记，检查产品的质量是否合格，几乎所有决策问题都可以用伯努利分布来描述。伯努利分布是最基本的一种分布。



伯努利分布的矩

$$E(X)=0\cdot(1-p)+1\cdot p=p$$

$$E(X^2)=0^2\cdot(1-p)+1^2\cdot p=p$$

$$Var(X)=E(X^2)-E(X)^2=p-p^2=p(1-p)$$



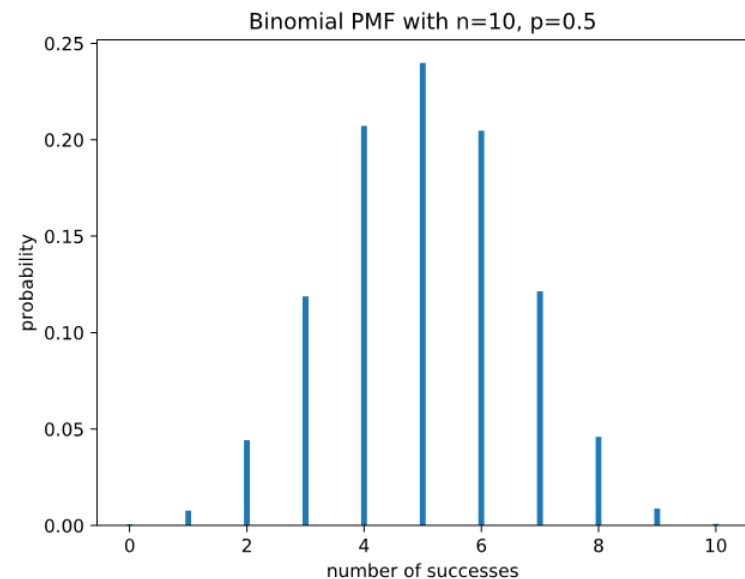
二项分布

二项分布 (Binomial Distribution) 是 n 个独立的伯努利试验的离散概率分布。

【二项分布】 设事件 A 在任一次试验中出现的概率为 p , 则在 n 重伯努利试验中事件 A 发生的次数 k 的取值为 $0, 1, \dots, n$ 。该随机变量 X 的概率分布为:

$$P(X=k) = C_n^k p^k (1-p)^{(n-k)}, k=0, 1, \dots, n$$

则称 X 服从参数为 n, p 的二项分布, 记为 $X \sim \mathbb{B}(n, p)$.





二项分布的矩

设 $X = \sum_{i=1}^n X_i$, 其中 X_i 服从伯努利分布

$$E(X) = E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i) = np$$

$$Var(X) = Var\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n Var(X_i) = np(1-p)$$



二项分布示例

【例】某人进行射击, 设每次射击的命中率为0.02, 独立射击400次, 试求至少击中两次的概率。

解: 将一次射击看成是一次试验. 设击中的次数为 X , 则 $X \sim \mathbb{B}(400, 0.02)$. X 的分布律为

$$P\{X=k\} = C_{400}^k (0.02)^k (0.98)^{(400-k)}$$

即得所求概率为

$$\begin{aligned} P\{X \geq 2\} &= 1 - P(X=0) - P(X=1) \\ &= 1 - (0.98)^{400} - 400(0.02)(0.98)^{399} = 0.9972 \end{aligned}$$

一个事件发生的概率不论多小, 只要不断重复试验下去, 事件迟早会出现的, 概率接近1。

二项分布示例

【例】设有80台同类型设备, 各台工作是相互独立的, 发生故障的概率都是0.01. 且一台设备的故障能由一个人处理. 考虑两种配备维修工人的方法, 其一是由4人维护, 每人负责20台; 其二是由3人共同维护80台. 试比较这两种方法在设备发生故障时不能及时维修的概率的大小.

解: 按第一种方法. 以 X 记“第1人维护的20台中同一时刻发生故障的台数”, 以 $A_i (i=1,2,3,4)$ 表示事件“第 i 人维护的20台中发生故障不能及时维修”, 则知80台中发生故障而不能及时维修的概率为

$$P(A_1 \cup A_2 \cup A_3 \cup A_4) \geq P(A_1) = P\{X \geq 2\}$$

而 $X \sim B(20, 0.01)$, 故有 $P(A_1 \cup A_2 \cup A_3 \cup A_4) \geq P\{X \geq 2\} = 1 - P\{X=0\} - P\{X=1\} = 0.0169$

按第二种方法. 以 Y 记80台中同一时刻发生故障的台数. 此时, $Y \sim B(80, 0.01)$, 故80台中发生故障而不能及时维修的概率为 $P\{Y \geq 4\} = 1 - \sum_{i=0}^3 C_{80}^i (0.01)^i (0.99)^{(80-i)} = 0.0087$

我们发现, 在后一种情况尽管任务重了(每人平均维护约27台), 但工作效率不仅没有降低, 反而提高了.

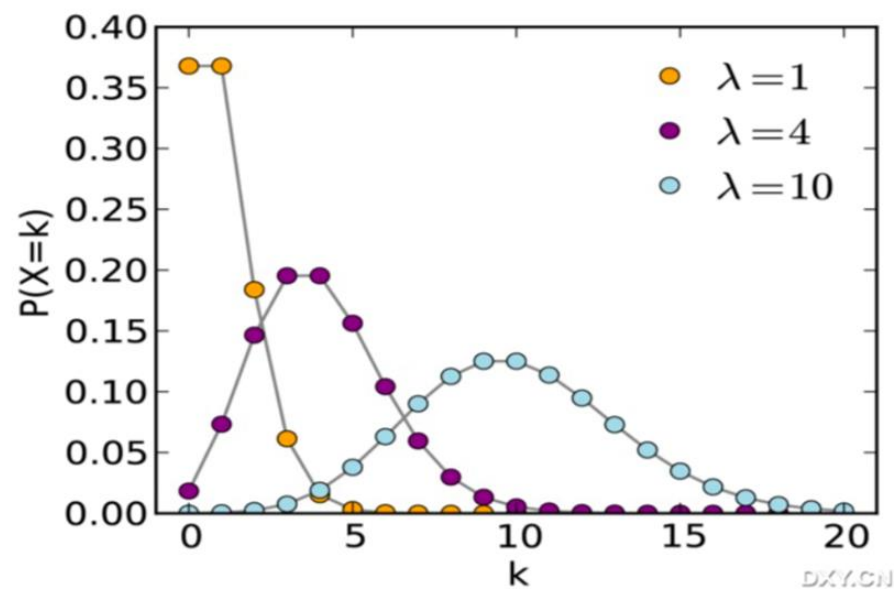
泊松分布

泊松分布 (Poisson Distribution) 是法国数学家泊松于1837年引入的。泊松分布适合于描述单位时间内随机事件发生的次数的概率分布。如某一服务设施在一定时间内受到的服务请求的次数, 系统出现的故障数、自然灾害发生的次数、DNA序列的变异数、放射性原子核的衰变数等等。

【泊松分布】 若随机变量 X 的概率分布为:

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, k = 0, 1, 2, \dots$$

则称 X 服从参数为 λ 的泊松分布, 记为 $X \sim \pi(\lambda)$.

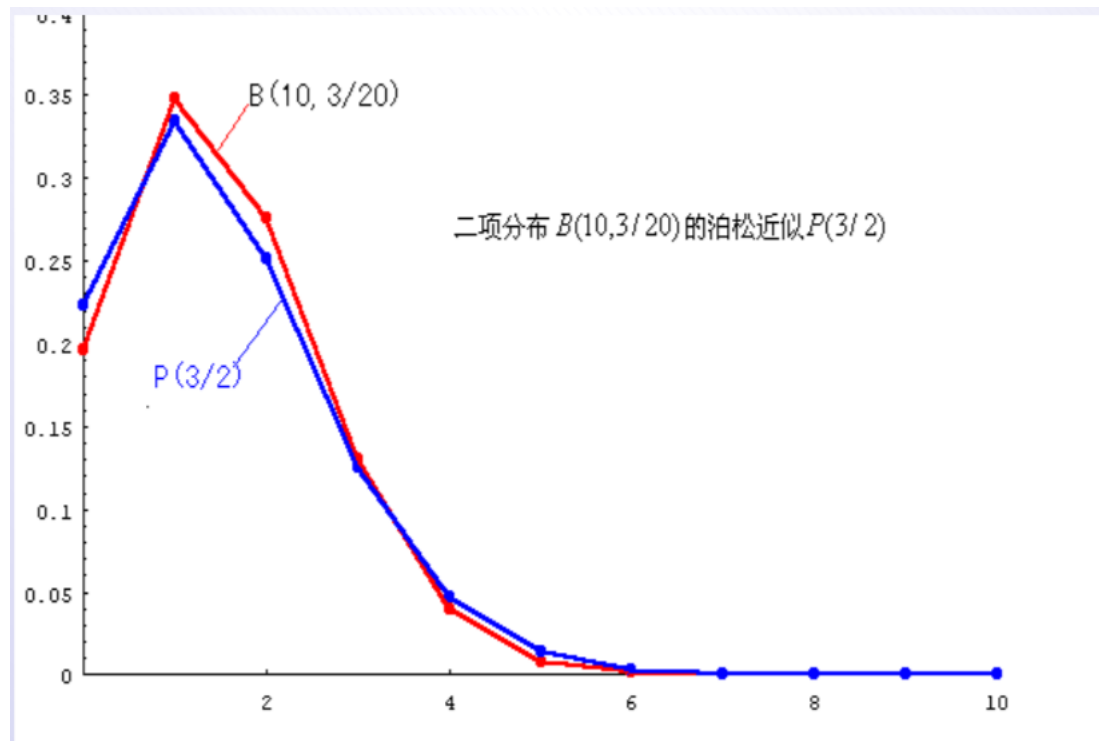


泊松逼近

在很多应用问题中, 我们常常有这样的伯努利试验, 其中, 相对地说, n 大, p 小, 而乘积 $\lambda = np$ 大小适中. 在这种情况下, 有一个便于使用的近似公式.

【泊松逼近】 在伯努利试验中, 以 p_n 代表事件 A 在试验中出现的概率, 如果 $np_n \rightarrow \lambda$, 则当 $n \rightarrow \infty$ 时,

$$\lim_{n \rightarrow \infty} \mathbb{B}(n, p_n) = \pi(\lambda)$$





泊松逼近证明

回顾二项分布公式: $P(X = k) = C_n^k p^k (1 - p)^{n-k}, p = \frac{\lambda}{n}$

证明:

$$\begin{aligned}\lim_{n \rightarrow \infty} P(X = k) &= \lim_{n \rightarrow \infty} \binom{n}{k} p^k (1 - p)^{n-k} \\&= \lim_{n \rightarrow \infty} \frac{n!}{(n-k)!k!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\&= \lim_{n \rightarrow \infty} \underbrace{\left[\frac{n!}{n^k (n-k)!} \right]}_F \underbrace{\left(\frac{\lambda^k}{k!}\right)}_{\rightarrow \exp(-\lambda)} \underbrace{\left(1 - \frac{\lambda}{n}\right)^n}_{\rightarrow \exp(-\lambda)} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-k}}_{\rightarrow 1} \\&= \lim_{n \rightarrow \infty} \underbrace{\left[\left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \dots \left(1 - \frac{k-1}{n}\right) \right]}_{\rightarrow 1} \underbrace{\left(\frac{\lambda^k}{k!}\right)}_{\rightarrow \exp(-\lambda)} \underbrace{\left(1 - \frac{\lambda}{n}\right)^n}_{\rightarrow \exp(-\lambda)} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-k}}_{\rightarrow 1} \\&= \left(\frac{\lambda^k}{k!}\right) \exp(-\lambda)\end{aligned}$$



泊松逼近计算

【例】假如生三胞胎的概率为 10^{-4} , 求在100000次生育中, 有0, 1, 2次生三胞胎的概率.

解:可看作伯努利试验; $n = 100000, p = 0.0001$, 所求的概率直接计算为

- $\mathbb{B}(0; 100000, 0.0001) = 0.000045378$
- $\mathbb{B}(1; 100000, 0.0001) = 0.00045382$
- $\mathbb{B}(2; 100000, 0.0001) = 0.0022693$

这时也可用泊松逼近, $\lambda = np = 10$, 而

- $\pi(0; 10) = 0.00004540$
- $\pi(1; 10) = 0.0004540$
- $\pi(2; 10) = 0.002270$



泊松逼近计算

【定理】 设 $X \sim \pi(\lambda)$, 则 $E(X) = \lambda$.

证:

$$\begin{aligned} E(X) &= \sum_{k=0}^{\infty} k \lambda^k k! e^{-\lambda} \\ &= \lambda^{-\lambda} \sum_{k=1}^{\infty} \lambda^{k-1} (k-1)! = \lambda e^{-\lambda} e^{\lambda} = \lambda \end{aligned}$$

【定理】 设 $X \sim \pi(\lambda)$, 则 $Var(X) = \lambda$.

证:

$$\begin{aligned} E(X^2) &= E[X(X-1) + X] = E[X(X-1)] + E(X) \\ &= \sum_{k=0}^{\infty} k(k-1) \lambda^k k! e^{-\lambda} + \lambda = \lambda^2 e^{-\lambda} \sum_{k=2}^{\infty} \frac{\lambda^{k-2}}{(k-2)!} + \lambda \\ &= \lambda^2 e^{-\lambda} e^{\lambda} + \lambda = \lambda^2 + \lambda \\ Var(X) &= E(X^2) - E(X)^2 = \lambda \end{aligned}$$

几何分布

【例】 从生产线上随机抽产品进行检测,设产品的次品率为 p , $0 < p < 1$, 若查到次品就停机检修,设停机时已检测到 X 只产品。 X 的概率分布律

设 A_i 为第 i 个抽到正品事件, A_i 相互独立, 则

$$P\{X = n\} = P\{A_1, \dots, A_{n-1}, \overline{A_n}\} = (1 - p)^{n-1}p$$

【几何分布】 如果随机变量 X 的概率分布为:

$$P(X = n) = (1 - p)^{n-1}p, n = 1, 2, \dots, 0 < p < 1$$

则称 X 服从几何分布, 记为 $X \sim \mathbb{G}(p)$



几何分布的矩

几何分布的数学期望($q = 1 - p$):

$$\begin{aligned} E(X) &= \sum_{i=1}^{\infty} i p q^{i-1} \\ &= 1p + 2pq + 3pq^2 + \dots + kpq^{k-1} + \dots \\ &= p(1 + 2q + 3q^2 + \dots + kq^{k-1} + \dots) \end{aligned}$$

令

$$\begin{aligned} S &= 1 + 2q + 3q^2 + \dots + kq^{k-1} + \dots \\ qS &= 1q + 2q^2 + 3q^3 + \dots + kq^k + \dots \end{aligned}$$

$$S - qS = (1 - q)S = 1 + q + q^2 + \dots + q^k + \dots = \frac{1}{1 - q}$$

$$\text{所以 } S = \frac{1}{(1-q)^2}, E(X) = pS = \frac{p}{(1-q)^2} = \frac{1}{p}$$



几何分布的矩

几何分布的方差($q = 1 - p$):

$$E(X^2) = \sum_{i=1}^{\infty} i^2 p q^{i-1}$$

$$= p(1 + 2^2 q + 3^2 q^2 + \dots + k^2 q^{k-1} + \dots)$$

令 $S = 1 + 2^2 q + 3^2 q^2 + \dots + k^2 q^{k-1} + \dots$, 对 $T' = S$, 即 T 关于 q 求一阶导得 S

$$T = q + 2q^2 + 3q^3 + \dots + kq^k + \dots = \frac{q}{(1-q)^2}$$

则 $S = T'$,

$$S = T' = \frac{(1-q)^2 + 2(1-q)q}{(1-q)^4} = \frac{1-q^2}{(1-q)^4} = \frac{1+q}{(1-q)^3}$$

$$\text{所以 } E(X^2) = pS = \frac{2-p}{p^2}$$

$$D(X) = E(X^2) - (EX)^2 = \frac{2-p}{p^2} - \frac{1}{p^2} = \frac{1-p}{p^2}$$

几何分布的无记忆性

【无记忆】 对于一个非负随机变量 X , 如果对于任意 $s, t \geq 0$ 有

$$P\{X > s + t | X > t\} = P\{X > s\}$$

则我们称 X 是无记忆的。

$$P\{X \leq n\} = 1 - q^n, P\{X > n\} = q^n$$

$$P\{X > n\} = \sum_{i=n+1}^{\infty} pq^{i-1} = pq^n \sum_{i=0}^{\infty} q^i = \frac{pq^n}{1-q} = q^n$$

因此无记忆性是几何分布所具有的一个有趣的性质。



超几何分布

一批产品共 N 件，含 M 件是次品，随机地从这 N 件产品中抽取 n 件产品，求恰有 k 件次品的概率。

【超几何分布】 如果随机变量 X 的概率分布为：

$$P(X = k) = \frac{C_M^k C_{N-M}^{n-k}}{C_N^n}, k = 1, 2, \dots, n$$

其中 N, M, n 均为正整数, 且 $M \leq N, n \leq N$, 则称 X 服从参数为 N, M, n 的超几何分布, 记为 $X \sim \mathbb{H}(N, M, n)$.



超几何分布近似

我们把二项分布与超几何分布作一比较。 N 件产品，有 M 件次品和 $N - M$ 件正品。如果每抽一件产品放回后，再抽下一件产品，如此有放回地随机地抽取 n 件，这是 n 重伯努利试验，那么所抽的 n 件产品的次品数 $X \sim \mathbb{B}(n, \frac{M}{N})$ 。当 $N \gg n$ 时，我们可以采用二项分布近似超几何分布。

超几何分布、二项分布和泊松分布都是重要的离散型随机变量的概率分布。有时，他们的概率计算会十分繁冗。当试验次数 n 很大时，可以推导出这三个分布间有一种近似关系式

$$\frac{C_M^k C_{N-M}^{n-k}}{C_N^n} \approx C_n^k p^k (1-p)^{n-k} \approx \frac{\lambda^k}{k!} e^{-\lambda}$$



超几何分布的矩

【超几何分布的数学期望】

$$n \frac{M}{N}$$

【超几何分布的方差】

$$n \frac{M}{N} \frac{N - M}{N} \frac{N - n}{N - 1}$$



总结

