

Stat 2911 Lecture Notes

Class 4 , 2017

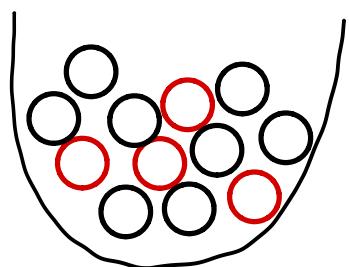
Uri Keich

© Uri Keich, The University of
Sydney

Class 4: Hypergeometric RV,
Fisher's Exact Test, joint PMF,
independent discrete RVs, deriving
the Poisson distribution

5) A hypergeometric RV X models the number of red balls in a sample of m balls drawn without replacement from an urn with r red balls and $n-r$ black balls.

A 3-parameter family (r, n, m)



$$\bar{X}(\Omega) \subset \{0, 1, \dots, r\}$$

All samples of size m are assumed equally likely:
equiprobable space.

For $k \in \{0, 1, \dots, r\}$,

$$P(\bar{X}=k) = \frac{\binom{r}{k} \binom{n-r}{m-k}}{\binom{n}{m}}$$

Fisher's Exact Test

(Rice 13.2)

Statistical Methods for Research Workers (1925)

30 convicted criminals with same sex twins, of which 13 are monozygotic (*identical*) twins and 17 are dizygotic twins.

10/13 monoz. twin siblings and 2/17 diz. twin siblings were themselves convicted.

Fisher: is there evidence of a genetic link?

twin's fate:

	convicted	not convicted	
mono	10	3	13
di	2	15	17
	12	18	

Assuming that, whether or not the sibling of a convicted criminal is also convicted does not depend on the type of the twin, we have a sample from a hypergeometric dist:

13 red balls (monoz. twin siblings)

17 black balls (diz. " "

We randomly sample 12 balls (convicted siblings).

$\bar{X} = \# \text{ of convicted monoz. siblings} \sim \text{HyperG}(13, 17, 12)$
 $(r) (n-r) (m)$

What is the prob. we will observe 10 or more red balls in the sample? (p-value)

$$\begin{aligned} P(\bar{X} \geq 10) &= \sum_{k=10}^{12} P(\bar{X}=k) \\ &= \sum_{k=10}^{12} \frac{\binom{13}{k} \binom{17}{12-k}}{\binom{30}{12}} \approx 0.000465 \end{aligned}$$

It seems very unlikely that there is no relation between the twin type and their conviction.

But did we establish that criminal mind is inherited?

- probabilistic statement
- conviction vs. truth "that face is up to no-good"
- sampling/ascertainment bias: finger print already in DB
- Identical twins feel tighter connections

Joint distribution

The joint pmf of the RVs X & y specifies their interaction:

$$p_{\bar{X}Y}(x, y) = P(\bar{X}=x, Y=y) \quad x, y \in \mathbb{R}$$

Note $\{\bar{X}=x, Y=y\} = \{\omega : \bar{X}(\omega)=x, Y(\omega)=y\}$

How do you know this is an event?

If X & Y are discrete RVs with

$$\{x_i\} = \bar{X}(\Omega) \text{ and } \{y_j\} = Y(\Omega)$$

then

$$p_{\bar{X}}(x) = P(\bar{X}=x) \stackrel{?}{=} \sum_j P(\bar{X}=x, Y=y_j) = \sum_j p_{\bar{X}Y}(x, y_j).$$

Similarly,

$$p_Y(y) = \sum_i p_{\bar{X}Y}(x_i, y).$$

$p_{\bar{X}}$ and p_Y are called the marginal pmfs.

Example (Rice 3.2)

A fair coin is tossed 3 times.

Let $X = \{ \text{HT on first toss} \} = \# \text{ of heads in first toss}$
 $Y = \text{total } \# \text{ of heads}$

$\Omega = \{ \text{HHH, HHT, ...} \}$ equiprobable space

$$|\Omega| =$$

x\y	0	1	2	3
0				
1				

$$\{ X=0, Y=0 \} = \{ \}$$

$$\Rightarrow P(X=0, Y=0) =$$

x\y	0	1	2	3
0	1/8	2/8	1/8	0
1	0	1/8	2/8	1/8
	1/8	3/8	3/8	1/8

$$X \sim$$

$$Y \sim$$

Are X and Y "independent"?

Independent RVs

The RV \underline{X} is ind. of y if knowing only the value of y does not give us any information about \underline{X} , in particular, it does not change the dist. of \underline{X} :

$$P(\underline{X} = \underline{x}_i | y = y_j) = P(\underline{X} = \underline{x}_i) \quad \forall i, j$$

Def.
 \Rightarrow

$$P(\underline{X} = \underline{x}_i, y = y_j) = P(\underline{X} = \underline{x}_i) P(y = y_j) \quad \forall i, j$$

$$\Leftrightarrow P_{\underline{X}y}(\underline{x}_i, y_j) = P_{\underline{X}}(\underline{x}_i) P_y(y_j) \quad \forall i, j$$

The joint prob factors into the product of the marginals.

More generally, the RVs $\underline{X}_1, \dots, \underline{X}_n$ are ind. if

$$P(\underline{X}_1 = \underline{x}_1, \dots, \underline{X}_n = \underline{x}_n) = P_{\underline{X}}(\underline{x}) = \prod_{i=1}^n P_{\underline{X}_i}(x_i) \quad \forall \underline{x} \in \mathbb{R}^n$$

$\underline{X} = (\underline{X}_1, \dots, \underline{X}_n)$ is a random vector and
 $\underline{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$.

Warning. Pairwise ind. of the RVs does not imply (mutual) ind.

The Poisson Dist.

Recall: $\bar{X} \sim \text{Poisson}(\lambda)$ if $P_{\bar{X}}(k) = e^{-\lambda} \frac{\lambda^k}{k!}$ $k \in \mathbb{N}$

A 1-parameter family: $\lambda > 0$

Where does it come from?

A Poisson dist. models the number of "events" that are registered in a certain time interval. For example, the number of

- particles emitted by a radioactive source in an hour
- incoming calls to a service center between 1-2 pm
- light bulbs burnt in a year
- fatalities from horse kicks in the Prussian cavalry

in one corps year (Bartkiewicz 1898, Rice p.215):

20 corps \times 10 years = 200 Corp years
 observed

# of deaths	count	frequency	Poisson approx.
0	109	0.545	0.543
1	65	0.325	0.331
2	22	0.110	0.101
3	3	0.015	0.021
4	1	0.005	0.003

What are the assumptions that lead to this part?

- 1) The dist. of the number of events at any time interval depends only on the interval's length/duration.
- 2) The # of events recorded in two disjoint time intervals are ind. of one another.
- 3) No two events are recorded at exactly the same time point.

Let $\bar{X}_{t,s} = \# \text{ of events in the time interval } (t, s]$ and let $\bar{X}_t = \bar{X}_{0,t}$.

Goal: find the dist. of \bar{X}_t .

We start with $P(\bar{X}_t=0)$ by studying $f(t) \stackrel{d}{=} P(\bar{X}_t=0)$.

$$\begin{aligned} f(t+s) &= P(\bar{X}_{t+s}=0) \\ &= P(\bar{X}_t=0, \bar{X}_{t,t+s}=0) \\ &\stackrel{?}{=} P(\bar{X}_t=0) P(\bar{X}_{t,t+s}=0) \\ &\stackrel{?}{=} P(\bar{X}_t=0) P(\bar{X}_{0,s}=0) = f(t) f(s) \end{aligned}$$

$$\Rightarrow f(t+s) = f(t) f(s) \quad \forall t, s \geq 0$$

What are the solutions of this functional equation?

One type is $f(t) = e^{\alpha t} \quad \alpha \in \mathbb{R}$

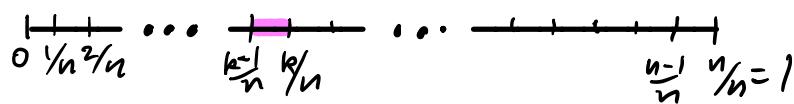
Other solutions are a "complete mess" and, in particular, they are unbounded, but $f(t) \in [0, 1] \quad \forall t \geq 0$ (why?).

For the same reason $\alpha < 0$ so $\alpha = -\lambda$ for $\lambda > 0$.

$$\Rightarrow P(X_t=0) = f(t) = e^{-\lambda t}$$

$$\text{In particular, for } t=1 \quad P(X_1=0) = e^{-\lambda} \quad (\overset{k}{e^{-\lambda}} \overset{k}{\cancel{\lambda^k / k!}}, k=0)$$

To find $P(X_i=j)$ for $j > 0$ consider dividing the time interval $(0, 1]$ to n intervals of size $\frac{1}{n}$:



Let $Y_n = \# \text{ of intervals } \left(\frac{k-1}{n}, \frac{k}{n}\right] \text{ in } (0, 1] \text{ in which an event occurred.}$

$$= \sum_{k=1}^n Y_n^k \quad \text{where} \quad Y_n^k = \mathbb{1}_{\{\bar{X}_{\frac{k-1}{n}, \frac{k}{n}} \geq 1\}}$$

Y_n^1, \dots, Y_n^n are iid Bernoulli (p_n) RVs where

$$\begin{aligned} p_n &= P(\bar{X}_{\frac{k-1}{n}, \frac{k}{n}} \geq 1) \\ &= P(\bar{X}_{\frac{1}{n}} \geq 1) \\ &= 1 - P(\bar{X}_{\frac{1}{n}} = 0) \\ &= 1 - e^{-\lambda n} \end{aligned}$$

$$\Rightarrow Y_n \sim \text{Binomial}(n, p_n = 1 - e^{-\lambda n})$$

Note: Generally, $Y_n \leq \bar{X}_i$ and $Y_n < \bar{X}_i$ if two events occur in the same interval.

However, $\lim_{n \rightarrow \infty} Y_n(\omega) = \bar{X}_i(\omega)$ $\forall \omega \in \Omega$ because no two events can occur at the same time.