

Stat 2911 Lecture Notes

Class 9 , 2017

Uri Keich

© Uri Keich, The University of
Sydney

Covariance, Variance of a sum of
RVs, Independence and 0
covariance, Variance of the
binomial, hypergeometric, and
negative binomial distributions
(Rice 4.3)

$$L' = \{ \underline{x} : \sum_i |x_i| p_x(x_i) < \infty \} \quad [\underline{x} \in L' \Leftrightarrow E|\underline{x}| < \infty]$$

L' is the vector space of RVs with finite mean.

$$L^2 = \{ \underline{x} : \sum_i x_i^2 p_x(x_i) < \infty \} \quad [\underline{x} \in L^2 \Leftrightarrow E\underline{x}^2 < \infty]$$

Generally, $L^2 \subsetneq L'$, where $\Leftrightarrow \underline{x} \in L'$

L^2 is the vector space of RVs with finite variance.

Generally, $\underline{x}, \underline{y} \in L' \not\Rightarrow \underline{x} \cdot \underline{y} \in L'$

Claim. Suppose $\underline{x}, \underline{y} \in L'$ are ind then $\underline{x} \cdot \underline{y} \in L'$
and $E(\underline{x}\underline{y}) = E(\underline{x})E(\underline{y})$

Claim. If $\underline{x}, \underline{y} \in L^2$ then $\underline{x} \cdot \underline{y} \in L'$

Claim. If $\underline{x} \in L^2$ $\alpha, \beta \in \mathbb{R}$ then

$$V(\alpha \underline{x} + \beta) = \alpha^2 \sigma^2$$

Want $V(\underline{x} + \underline{y})$

Def. For jointly distributed RVs $\bar{x}, y \in L'$ we define the covariance of \bar{x} and y as

$$\text{Cov}(\bar{x}, y) = E[(\bar{x} - \mu_x)(y - \mu_y)],$$

provided $(\bar{x} - \mu_x)(y - \mu_y) \in L'$ and where

$$\mu_x = E(\bar{x}), \quad \mu_y = E(y).$$

Note the symmetry: $\text{Cov}(\bar{x}, y) = \text{Cov}(y, \bar{x})$.

When does the covariance exist?

Claim. If $\bar{x}, y, \bar{x} \cdot y \in L'$ then $\text{Cov}(\bar{x}, y)$ exists and

$$\text{Cov}(\bar{x}, y) = E(\bar{x}y) - E(\bar{x})E(y)$$

Cor. If $\bar{x}, y \in L^2$ then the conclusion of the claim holds.

Proof. $(\bar{x} - \mu_x)(y - \mu_y) = \bar{x}y - \mu_x y - \mu_y \bar{x} + \mu_x \mu_y$

Since $\bar{x}y, \bar{x}, y, 1$ (constant) $\in L'$ the R.H.S $\in L'$

$\Rightarrow (\bar{x} - \mu_x)(y - \mu_y) \in L'$ and its expected value is

$$\begin{aligned} \text{Cov}(\bar{x}, y) &= E(\bar{x}y) - \mu_x E(y) - \mu_y E(\bar{x}) + \mu_x \mu_y \\ &= E(\bar{x}y) - \mu_x \mu_y \end{aligned}$$

□

Cor. If $X, Y \in L'$ are ind. then

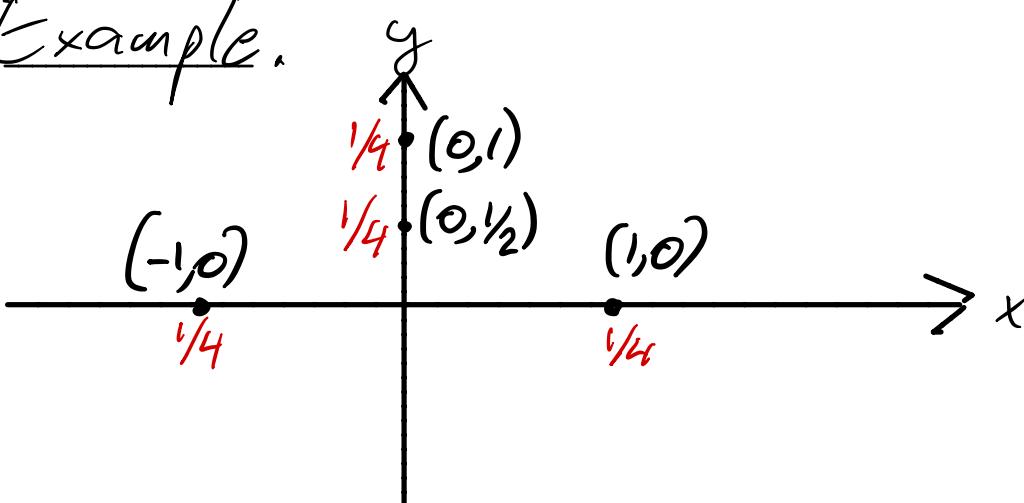
$$\text{Cov}(X, Y) =$$

Or, independent RVs are uncorrelated
(having 0 covariance)

What about the converse?

Covariance is a measurement of the linear dependence between RVs, which is substantially more restricted than probabilistic dependence.

Example.



Can't predict
y using a linear
function of x
 \nRightarrow can't
predict y
from X

The joint pmf p_{xy} assigns a mass of $\frac{1}{4}$ for the four depicted points (and 0 elsewhere)

Check: $\text{Cov}(X, Y) = 0$ but X and Y are not ind.

$E(\bar{x}+y) = E(\bar{x}) + E(y)$, is $V(\bar{x}+y) = V(\bar{x}) + V(y)$?

Clearly we need $\bar{x}, y \in L^2$ for the RHS.

$\Rightarrow \bar{x}+y \in L^2$ so $V(\bar{x}+y)$ exists and is finite.

If $\bar{x} = -y$ then $V(\bar{x}+y) = V(\bar{x}-\bar{x}) =$

$V(\bar{x}+y)$ can be smaller than $V(\bar{x})+V(y)$

If $\bar{x} = y$ then $V(\bar{x}+y) = V(2\bar{x}) =$

$V(\bar{x}+y)$ can be larger than $V(\bar{x})+V(y)$

Claim. If $\bar{x}, y \in L^2$ then

$$V(\bar{x}+y) = V(\bar{x}) + V(y) + 2\text{Cov}(\bar{x}, y)$$

Proof. Note that $\text{Cov}(\bar{x}, y)$ is well defined.

$$E(\bar{x}+y)^2 = E(\bar{x}^2 + 2\bar{x}y + y^2) = E(\bar{x}^2) + 2E(\bar{x}y) + E(y^2)$$

$$E^2(\bar{x}+y) = [E(\bar{x}) + E(y)]^2 = E^2(\bar{x}) + 2E(\bar{x})E(y) + E^2(y)$$

subtract

$$V(\bar{x}+y) = V(\bar{x}) + V(y) + 2\text{Cov}(\bar{x}, y)$$

$$= V(\bar{x}) + V(y) + \text{Cov}(\bar{x}, y) + \text{Cov}(y, \bar{x})$$

By induction, if $\bar{X}_1, \dots, \bar{X}_n \in L^2$ then

$$V\left(\sum_1^n \bar{X}_i\right) = \sum_1^n V(\bar{X}_i) + \sum_{i \neq j} \text{Cov}(\bar{X}_i, \bar{X}_j)$$

For the proof you will need the linearity:

$$\text{Cov}(\bar{X} + \bar{Y}, \bar{Z}) = \text{Cov}(\bar{X}, \bar{Z}) + \text{Cov}(\bar{Y}, \bar{Z})$$

Cor. If \bar{X}_i are ind. L^2 -RVs then

$$V\left(\sum_1^n \bar{X}_i\right) = \sum_1^n V(\bar{X}_i)$$

Examples.

1) \bar{X}_i are iid Bernoulli(p) RVs

Let $S_n^1 = \sum_1^n \bar{X}_i$ then

$$\begin{aligned} V(S_n^1) &= \sum_1^n V(\bar{X}_i) \\ &= \sum_1^n p(1-p) \\ &= n p (1-p) \end{aligned}$$

Note that $S_n \sim$

2) Let $\bar{X} = \#$ of red balls in a sample of m balls from an urn with r red balls and $n-r$ black balls. Then $\bar{X} \sim \text{hyper}(n, r, m)$. What are $E(\bar{X})$, $V(\bar{X})$?

We can compute those directly from the definitions, using the pmf. Alternatively:

For $i=1, 2, \dots, m$, let $X_i = 1_{\text{if } i^{\text{th}} \text{ drawn ball is red}}$.

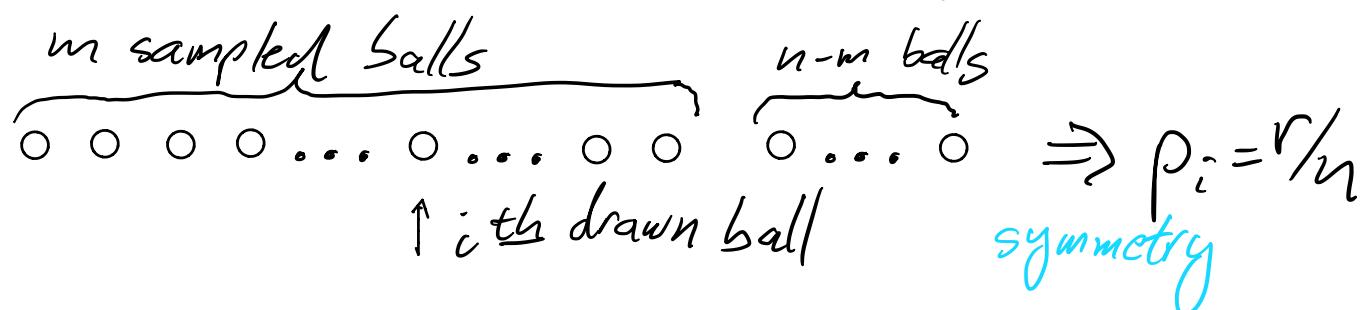
Clearly, $X_i \sim \text{Bernoulli}(p_i)$, where

$$p_1 =$$

$$\begin{aligned} p_2 &= P(\bar{X}_2 = 1 | \bar{X}_1 = 1) p(\bar{X}_1 = 1) + P(\bar{X}_2 = 1 | \bar{X}_1 = 0) p(\bar{X}_1 = 0) \\ &= \frac{r-1}{n-1} \cdot \frac{r}{n} + \frac{r}{n-1} \cdot \frac{n-r}{n} = \dots = \end{aligned}$$

$$p_3 = \dots = \text{why?}$$

The sampling is consistent with randomly sorting the n balls in a row and picking the first m balls.



Clearly, $\bar{X} =$

$$\Rightarrow E(\bar{X}) = \sum_1^m E(X_i) =$$

Are \bar{X}_i ind.?

The \bar{X}_i are not ind. : $\bar{X}_i = 1$ implies that the proportion of remaining red balls has decreased

$$V(\bar{X}) = \sum_1^m V(\bar{X}_i) + \sum_{i \neq j} \text{Cov}(\bar{X}_i, \bar{X}_j)$$

$$V(\bar{X}_i) = r/n (1 - r/n)$$

$$\text{Cov}(\bar{X}_i, \bar{X}_j) = E(\bar{X}_i \bar{X}_j) - \underbrace{E(\bar{X}_i) E(\bar{X}_j)}_{(r/n)^2}$$

$$E(\bar{X}_i \bar{X}_j) = 1 \cdot P(\bar{X}_i = 1, \bar{X}_j = 1)$$

$$= \frac{r(r-1)}{n(n-1)} \quad \text{if } i \neq j$$

$$\begin{aligned} & \text{if } i \neq j \\ \Rightarrow \text{Cov}(\bar{X}_i, \bar{X}_j) &= \frac{r(r-1)}{n(n-1)} - \frac{r^2}{n^2} = \frac{r}{n} \left(\frac{r-1}{n-1} - \frac{r}{n} \right) \\ &= -\frac{r}{n} \left(1 - \frac{r}{n} \right) \frac{1}{n-1} \quad \begin{array}{l} \text{(does the sign make sense?)} \\ \text{make sense?} \end{array} \end{aligned}$$

$$\begin{aligned} \Rightarrow V(\bar{X}) &= m \frac{r}{n} \left(1 - \frac{r}{n} \right) + m(m-1) \left(-\frac{r}{n} \right) \left(1 - \frac{r}{n} \right) \frac{1}{n-1} \\ &= m \frac{r}{n} \left(1 - \frac{r}{n} \right) \left(1 - \frac{m-1}{n-1} \right) \end{aligned}$$

It is instructive to compare this expression to sampling with replacement.

Let $Y = \#$ of red balls in a sample of m balls taken with replacement where $m \leq n$.

Clearly, $Y \sim$

Therefore,

$$E(Y) = E(\bar{X})$$

$$V(Y) = m \frac{r_n}{n} \left(1 - \frac{r_n}{n}\right) \left(1 - \frac{m-1}{n-1}\right) = V(\bar{X})$$

and $>$ if $m > 1$

Does it make sense?

$$m=1: V(\bar{X}) = V(Y) = \frac{r_n}{n} \left(1 - \frac{r_n}{n}\right)$$

as both \bar{X} and Y are Bernoulli(r/n) RVs

$$m=n: V(\bar{X}) =$$

Fix m and let $n \rightarrow \infty$ s.t. $r_n/n \rightarrow p$, then

$$V(\bar{X}_n) = m \frac{r_n}{n} \left(1 - \frac{r_n}{n}\right) \left(1 - \frac{m-1}{n-1}\right) \xrightarrow{n} mp(1-p)$$

$$V(Y_n) = m \frac{r_n}{n} \left(1 - \frac{r_n}{n}\right) \xrightarrow{n} mp(1-p)$$

Makes sense?

3) $N \sim \text{negative binomial}(r, p)$ [NB(r, p)]

Claim. If \bar{X}_i are iid geometric(p) RVs then

$$\sum_{i=1}^r \bar{X}_i \sim NB(r, p).$$

Proof. Algebraic, using convolutions (problem set).

Alternatively, consider \bar{X}_i as counting the number of additional trials to the i th success.

By definition,

$$N = \sum_{i=1}^r \bar{X}_i \sim NB(r, p).$$

At the same time it is intuitively clear that the \bar{X}_i are iid geometric(p) RVs. □

It follows that

$$\begin{aligned} E(N) &= \sum_{i=1}^r E(\bar{X}_i) \\ &= r/p. \end{aligned}$$

$$\begin{aligned} V(N) &= \sum_{i=1}^r V(\bar{X}_i) \\ &\approx \frac{r(1-p)}{p^2}. \end{aligned}$$

Alternatively, compute using the $NB(r, p)$ pdf.