

## COVARIANCE AND CORRELATION

Rice P. 129  
 Lindgren P. 134  
 BAIN AND ENGELEN P. 179  
 188

If we consider the 2 diml r. var  $X, Y$ , we can introduce the concept of a MIXED MOMENT (or more usually — a PRODUCT MOMENT) :

$$E[(X - a)^r (Y - b)^s]$$

for  $r, s$  integers  $> 0$

We shall be mainly concerned with one particular product moment :

$$E[(X - \mu_x)(Y - \mu_y)]$$

— known as the covariance of  $X$  and  $Y$  — or the covariance of the Bivariate distribution.

— equivalent form:  $E[XY] - (EX)(EY)$

NOTATION :  $\sigma_{xy}$  or  $\text{Cov}(X, Y)$

The covariance provides some measure of the extent to which  $X$  and  $Y$  co-vary.

NOTICE that if  $X, Y$  are indept, then  $\text{Cov}(X, Y) = 0$  (but not the converse)

## CORRELATION COEFF.

However the covariance is a poor measure of covariation because it is scale-dependent

e.g. if rv  $X, Y$  took values in inches the  $\text{Cov}(X, Y)$  would be 144 times what it would be if  $X, Y$  were measured in feet.

To remove this defect, the CORR COEFF is introduced:

DEFN :

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

NOTE

In attempting to get an intuitive understanding of what the corr. coeff measures, it is important to realize that it measures a PARTICULAR KIND OF DEPENDENCE

— but there are other kinds of dependence which are not measured by it.

(We have already pointed out that it is possible to have  $\sigma_{xy} = 0$  when  $x, y$  are DEPENDENT.)

SOME RESULTS ON  $\rho$ 

$$\text{II} \quad |\rho_{xy}| \leq 1$$

Proof : Suppose  $U, V$  are any random vars

Consider  $(U - kV)^2 \rightarrow$  a nonnegative r.var  
( $k$  real const)

$$E[U - kV]^2 = E[U^2] - 2k E[UV] + k^2 E[V^2]$$

and this is  $\geq 0$

Because this is  $\geq 0$ , we see that we must have

$$\begin{cases} A = E[V^2] \\ B = -2E[UV] \\ C = E[U^2] \end{cases}$$

$$\text{i.e. } [E(UV)]^2 \leq E[U^2]E[V^2]$$



Now choose  $U = X - EX$

$$V = Y - EY$$

so that  $\{E[(X - EX)(Y - EY)]\}^2 \leq \sigma_x^2 \sigma_y^2$

i.e.  $|\text{Cov}(x, y)| \leq \sigma_x \sigma_y$

and so  $|\rho_{xy}| \leq 1$

### RESULT 2

$$\text{if } \rho^2 = 1,$$

Then  $X$  and  $Y$  are linearly related  
with probability 1  
(i.e.  $\text{Prob}[(X - E(X)) = k_0(Y - EY)] = 1$ )

Proof:

$E[UV]^2 = [E(U)]^2 [E(V^2)]$  if  $\rho^2 = 1$ , then the function  $E[\bar{U} - k\bar{V}]^2$

$\Rightarrow B^2 = 4AC$  has the following appearance:

i.e.  $\exists$  some value of  $k$

for which

$$E[\bar{U} - k_0\bar{V}]^2 = 0 \quad (B^2 - 4AC = 0)$$

Now, we choose  $\frac{\bar{U}}{\bar{V}} = X - E(X)$

$$\therefore \frac{\bar{U}}{\bar{V}} = Y - E(Y)$$

$$\therefore \text{that } E\bar{U} = E\bar{V} = E(\bar{U} - k\bar{V}) = 0$$

$$\text{And so } \text{Var}[\bar{U} - k_0\bar{V}] = 0$$

But, arising out of Chebyshev's Ineq, we saw that  
if  $\text{Var } X = 0$ , then  $(X = \mu)$  will prob. 1

In this case

$$(X - E(X)) - k_0(Y - E(Y)) = 0 \quad \text{with prob. 1}$$

$$\text{i.e. } X - E(X) = k_0(Y - E(Y)) \quad \text{with prob. 1}$$

i.e.  $X$  and  $Y$  are linearly related with prob. 1

## The Linear Relationship

Let's say  $k_0$  what is  $k_0 = ?$  : Just solve the Quadratic Eqn

$$\begin{aligned} \text{In fact } k_0 &= -\frac{B}{2A} \quad (B^2 - 4AC = 0) \\ &= \frac{2E[(X - \bar{X})(Y - \bar{Y})]}{2E(Y - \bar{Y})^2} \\ &= \frac{\sigma_{XY}}{\sigma_Y^2} ; (\text{Since } \rho^2 = 1) \text{ thus } \Rightarrow \end{aligned}$$

Thus the linear return is  $\frac{\sigma_{XY}}{\sigma_Y^2} = \frac{\sigma_X}{\sigma_Y}$

$$X - \bar{X} = \frac{\sigma_{XY}}{\sigma_Y^2} (Y - \bar{Y})$$

$$\left( \text{Since } \rho^2 = 1, \sigma_{XY} = \pm \sigma_X \sigma_Y \right)$$

$$= \pm \frac{\sigma_X}{\sigma_Y} (\bar{Y} - \bar{Y})$$

## SLOPE OF THE LINE

If  $\rho = 1$ , then  $\sigma_{XY} > 0$

and the slope of the line  $(\frac{\sigma_{XY}}{\sigma_Y^2})$  is +ve

If  $\rho = -1$  ( $\because \sigma_{XY} < 0$ )

The slope of the line will be -ve.

$$\sigma_{XY} = -\sigma_X \sigma_Y$$

NOTE : If  $Y = AX + B$  ( $A, B$  const.)

Then  $\rho = 1$

and  $\rho = +1$  if  $A > 0$

$\rho = -1$  if  $A < 0$

$\rho = 0$  DOES NOT  $\Rightarrow$  INDEPENDENCE

Counterexample : 3 coin tosses

$\bar{X} = \# \text{ heads}$ ;  $\# \text{ runs} = \bar{Y}$

$\bar{Y}$	1	2	3
0	$\frac{1}{8}$	0	0
1	0	$\frac{3}{8}$	$\frac{1}{8}$
2	0	$\frac{3}{8}$	$\frac{1}{8}$
3	$\frac{1}{8}$	0	0

$\frac{1}{4}$
$\frac{1}{2}$
$\frac{1}{4}$
$\frac{1}{8}$

$$f(x,y) \neq f_{\bar{X}}(x) f_{\bar{Y}}(y)$$

$$\text{Cov}(\bar{X}, \bar{Y}) = E[(\bar{X} - \mu_{\bar{X}})(\bar{Y} - \mu_{\bar{Y}})]$$

$$\mu_{\bar{X}} = \frac{3}{2}$$

$$\mu_{\bar{Y}} = 2$$


$$\begin{aligned} E(\bar{X}\bar{Y}) &= 0 \left( \frac{1}{8} \right) \\ &\quad + 2 \left( \frac{3}{8} \right) + 3 \left( \frac{1}{8} \right) \\ &\quad + 4 \left( \frac{3}{8} \right) + 6 \left( \frac{1}{8} \right) \\ &\quad + 3 \left( \frac{1}{8} \right) \\ &= \frac{24}{8} = 3 \end{aligned}$$

$$\begin{aligned} \text{Thus } \text{Cov}(\bar{X}, \bar{Y}) &= 3 - \left( \frac{3}{2} \right)^2 \\ &= 0 \end{aligned}$$

NOTE :  $\rho = 0 \Rightarrow$  linear covariation is absent

However this does not mean that  
 $\bar{X}$  and  $\bar{Y}$  are <sup>independently</sup> unrelated

There are other kinds of dependence apart from  
 Linear dependence — but these other kinds  
 may not be detected by  $f$

## OTHER CHARACTERISTICS OF PROB. DISTRIBUTIONS

We have pointed out that the Mean Value provides us with a AVERAGE — or measure of location for the random variable. (CENTRAL VALUE)

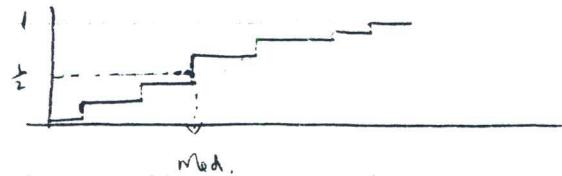
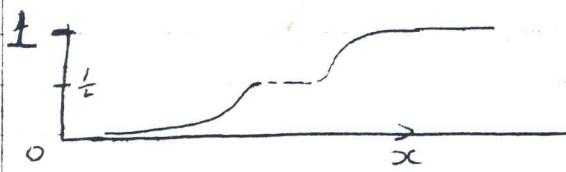
However there are other quantities which are also used to indicate the location of a prob. distrib.

1) Median : That point  $x$  for which

$$P[X \leq x] = \frac{1}{2} \quad F(x) = \frac{1}{2}$$

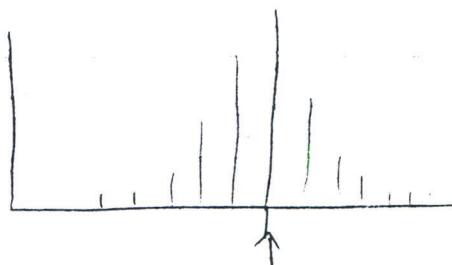
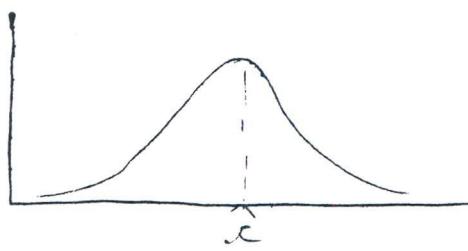
and  $1 - F(x) = \frac{1}{2}$

NOTE : If  $F(x)$  is never const on an interval, the median is uniquely defined



2) Mode

The value of  $x$  for which  $f(x)$  is maximum (either for contin. or discrete).



## OTHER MEASURES OF DISPERSION : QUANTILES

We can define quantiles  $x_\alpha$  as follows

$$F(x_\alpha) = \alpha \quad \text{for } 0 < \alpha < 1$$

Usual  $\alpha$  values are :  $0.01, 0.02, \dots, 0.99 \Rightarrow$  PERCENTILES

$0.1, 0.2, \dots, 0.9 \Rightarrow$  DECILES

$0.25, 0.5, 0.75 \Rightarrow$  QUARTILES

Dispersion Measures :  $P_{90} - P_{10}$  (10-90 PERCENTILE RANGE)  
 $Q_3 - Q_1$  (Interquartile range)

## SOME HIGHER MOMENTS

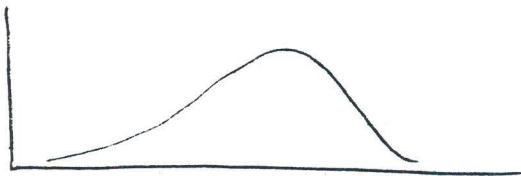
:  $\mu_3 + \mu_4$

SKEWNESS  
& KURTOSIS

A few words about  $\mu_3$  and  $\mu_4$

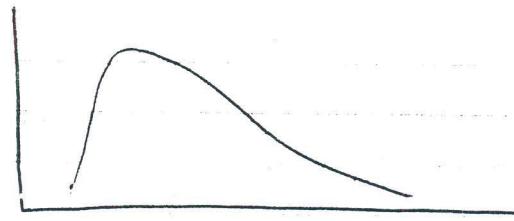
$\mu_3$  is often termed a measure of SKEWNESS (or ASYMMETRY)

It can be shown that symmetric distrib. will have  $\mu_3 = 0$  and that pdfs such as :



have negative  $\mu_3$

and



have +ve  $\mu_3$

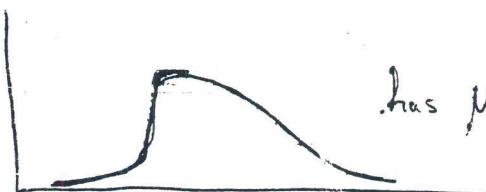
MOOD  
& GR.

P. 109-110

(In fact  $\frac{\mu_3}{\sigma^3}$  is termed the Coeff. of SKEWNESS)

However, a knowledge of  $\mu_3$  gives little information as to the shape of the distrib.:

e.g.



has  $\mu_3 = 0$

$\mu_4$

: This is regarded as a Measure of KURTOSIS  
— i.e. measures the extent of PEAKEDNESS in the distrib.

$$\text{Coeff. of KURTOSIS} = \frac{\text{Index of}}{\mu_4} = \frac{\mu_4}{\sigma^4}$$

$$\gamma_2 = \frac{\mu_4}{\sigma^4} - 3$$

For Normal Distrib.,  $\gamma_2 = 0$

Excess of  
Kurtosis

IF  $V(X) = 0$ , THEN  $X = \mu$  WITH PROBABILITY 1

Proof :

PAGE 7a

$$\text{Using } P[g(X) \geq c] \leq \frac{E[g(X)]}{c} \quad \text{for } c > 0$$

$$\text{choose } g(x) = (x - \mu)^2 \quad \text{and } c = t^2 \quad \text{and } g(x) \geq$$

$$\Rightarrow P[|X - \mu| \geq t] \leq \frac{\sigma^2}{t^2} \quad (*)$$

Now, Suppose  $P(X = \mu) < 1$

WE WILL TRY TO  
GET A CONTRADICTION

then there must be probability associated  
with values other than  $\mu$ ,

so that for some  $\varepsilon > 0$ ,

$$P[|X - \mu| \geq \varepsilon] > 0 \quad (*)$$

$$\text{But } P[|X - \mu| \geq \varepsilon] \leq \frac{\sigma^2}{\varepsilon^2} \quad (\text{from * above})$$

and using the fact that  $\sigma^2 = 0$ , we have

$$P[|X - \mu| \geq \varepsilon] = 0$$

and this contradicts (\*) above.

So, we must have

$$P(X = \mu) \text{ with probability 1.}$$

## VARIANCE OF SUMS OF RANDOM VARS

Lindgren P 137  
740

Let  $X_1, X_2, \dots, X_n$  be  $n$  random vars  
and let  $Z = \sum_{i=1}^n a_i X_i$  ( $a_i$ , real const.)

Then we want an expression for  $\text{Var}(Z)$

$$\begin{aligned} \text{Now } E(Z) &= E\left(\sum a_i X_i\right) \\ &= \sum_{i=1}^n a_i E X_i \quad (\text{Props of } E) \end{aligned}$$

$$\begin{aligned} \text{V}(Z) &= E[Z - E Z]^2 \\ &= E\left[\sum a_i X_i - \sum a_i (E X_i)\right]^2 \\ &= E\left[\sum a_i (X_i - E X_i)\right]^2 \\ &= E\left[\sum_{i=1}^n a_i^2 (X_i - E X_i)^2 + \right. \\ &\quad \left. + 2 \sum_i \sum_{j < i} a_i a_j (X_i - E X_i)(X_j - E X_j)\right] \\ &= \sum_{i=1}^n a_i^2 E(X_i - \mu_{X_i})^2 + 2 \sum_i \sum_{j < i} a_i a_j E(X_i - E X_i)(X_j - E X_j) \\ &= \sum_{i=1}^n a_i^2 \text{Var } X_i + 2 \sum_i \sum_{j < i} a_i a_j \text{Cov}(X_i, X_j) \end{aligned}$$

Another useful result:

$$Z = \sum_{i=1}^n a_i X_i ; \bar{W} = \sum_{i=1}^n b_i X_i$$

$$\text{Cov}(\bar{Z}, \bar{W}) = \sum_{i=1}^n a_i b_i \text{Var } X_i + 2 \sum_i \sum_{j < i} a_i b_j \text{Cov}(X_i, X_j)$$

$$\text{Proof: } \text{Cov}(\bar{Z}, \bar{W}) = E[(\bar{Z} - E \bar{Z})(\bar{W} - E \bar{W})]$$

$$= E\left[\left\{\sum a_i (X_i - \mu_i)\right\} \left\{\sum b_j (X_j - \mu_j)\right\}\right]$$

$$= E\left[\sum a_i b_i (X_i - \mu_i)^2 + 2 \sum_i \sum_{j < i} a_i b_j (X_i - \mu_i)(X_j - \mu_j)\right]$$

$$= \sum_{i=1}^n a_i b_i \text{Var}(X_i) + 2 \sum_i \sum_{j < i} a_i b_j \text{Cov}(X_i, X_j)$$

### Special cases

$$\underline{a_i = 1 \forall i} \quad \text{Var}(\sum \bar{X}_i) = \sum \text{Var}(\bar{X}_i) + 2 \sum_i \sum_{i < j} \text{Cov}(\bar{X}_i, \bar{X}_j)$$

If  $\text{Cov}(\bar{X}_i, \bar{X}_j) = 0$  for all pairs  $\bar{X}_i, \bar{X}_j$ .

Then  $\text{Var}(\sum \bar{X}_i) = \sum \text{Var } \bar{X}_i$

$$\text{Var}\left(\sum_{i=1}^n (a_i \bar{X}_i)\right) = \sum_{i=1}^n a_i^2 \text{Var } \bar{X}_i$$

$$\text{Cov}(\bar{Z}, \bar{W}) = \sum a_i b_i \text{Var } \bar{X}_i$$

Stronger  
 Cond'n;  
 all  $\bar{X}_i$  indept

(MEYER 148-150)  
 L. 114-115  
 BAIN/ENGELHARDT 190-194

## CONDITIONAL EXPECTATION

Suppose  $X, Y$  are r.v.s with pdf  $f(x, y)$

and  $f(x|y)$  is condit. PDF for  $X$  for given  $y$

Then :

DEFN We define the COND. EXPECTATION of  $X$  given  $y = y$

as  $E(X|y) = \int_{-\infty}^{\infty} x f(x|y) dx$  CONTINUOUS CASE

and

$$E(X|y) = \sum_{x_i} x_i f(x_i|y) \quad \text{DISCRETE CASE}$$

$E(X|y)$  is a fn of  $y \rightarrow$  say  $u(y)$ . We can then consider  $u(y)$

or  $E(u(Y))$

Useful Result :  $E_y [ E_x [ X | Y ] ] = E[X]$

Proof :

$$E_x [ X | Y ] = \int_{-\infty}^{\infty} x f(x|y) dx$$

$$= \underbrace{\int_{-\infty}^{\infty} x \frac{f(x,y)}{f_y(y)} dx}_{\text{This is a function of } y} \quad \text{not a random var.}$$

$\rightarrow u(y)$

Now let us consider  $u(Y)$ , which is what we mean by  $E(X|Y)$ ,  
 — and find  $E_y (u(Y))$

$$E_y [ u(Y) ] = E_y [ E_x [ X | Y ] ] = \int_{-\infty}^{\infty} E(x|y) f_y(y) dy$$

$$= \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} x \frac{f(x,y)}{f_y(y)} dx \right] f_y(y) dy \quad \text{stat}$$

Reversing the order of integration (permissible if the expectations exist.)

$$= \iint_{-\infty}^{\infty} x f(x,y) dy dx = \int x f(x) dx$$

$$= E(X)$$

Exa

Shipments of tomatoes arrive in a shop each day. However the number of tomatoes in the shipment varies and letting  $N$  denote the number of tomatoes, we shall assume that the distribution of  $N$  is as follows.

$n$	25	26	27	28	29	30	31	32
Prob( $N=n$ )	0.05	0.1	0.1	0.2	0.3	0.15	0.05	0.05

The shopkeeper knows from experience that the average proportion of bad tomatoes is 10%.

What is the expected number of bad tomatoes per day?

$$\frac{\bar{X}}{N} = \# \text{ bad tomatoes per day}$$

$$N = \# \text{ tomatoes in shipment}$$

$$\text{We want } E(\bar{X})$$

$$\text{We use } E(\bar{X}) = E(E(\bar{X}|N))$$

$$E(\bar{X}|n) = 0.10n$$

$$\begin{aligned} E_n[E(\bar{X}|N)] &= E_n[0.10N] \\ &= 0.10 E(N) \end{aligned}$$

$$\text{now } E(N) = 28.5$$

$$\text{Thus } E(\bar{X}) = 2.85$$

1.25
2.60
2.70
5.60
8.70
4.50
1.55
1.60
<hr/> 28.50

NOTE For any real valued function  $g(x, y)$ , it can be shown (using the same procedure as above) that

$$\begin{aligned} E[g(\bar{X}, Y)] &= E_Y \left[ E_{\bar{X}|Y}(g(\bar{X}, Y)|Y) \right] \\ &= E_{\bar{X}} \left[ E_{Y|\bar{X}}(g(\bar{X}, Y)|\bar{X}) \right] \end{aligned}$$

$$E[g(\bar{X}, Y)] = \iint g(x, y) f(x, y) dy dx$$

$$= \iint g(x, y) \frac{f(x, y)}{f_x(x)} f_x(x) dy dx$$

$$= \int \left[ \iint g(x, y) f(y|x) dy \right] f_x(x) dx$$

$$E_x \left[ E[g(\bar{X}, Y)|x] \right]$$

CONDITIONAL VARIANCELindgren p 130 $X, Y$  has joint distib.  $f(x, y)$ Defn

$$\text{Var}(X|y) = E[(X - \mu_{x|y})^2 | y]$$

$$\mu_{x|y} = E(X|y)$$

This is again a function of  $y$ , say  $u(y)$ Thus we can consider  $u(y)$ , written as  $\text{Var}(X|Y)$ We might expect that  $E_y[\text{Var}(X|Y)] = \text{Var}(X)$  — but Noto investigate this:

We need the following result:

$$E[(X - a)^2] = \text{Var } X + (\mu - a)^2$$

PARALLEL  
AXIS  
THEOREM

Now

$$E[(X - \mu_x)^2 | y] = \text{Var}(X|y) + (\mu_{x|y} - \mu_x)^2$$

i.e.

$$\int_{-\infty}^{\infty} (x - \mu_x)^2 \frac{f(x, y)}{f_y(y)} dx =$$

$\underbrace{g(y)}$  say

Let us take the expected value of  $g(y)$ 

$$\Rightarrow \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} (x - \mu_x)^2 \frac{f(x, y) dx}{f_y(y)} \right] f_y(y) dy = E_y \text{Var}(X|y) + E_y [\mu_{x|y} - \mu_x]^2$$

$$\underbrace{\iint (x - \mu_x)^2 f(x, y) dx dy} =$$

$$\text{Var}(X) = E_y \text{Var}(X|y) + \text{Var}(\mu_{x|y})$$

$$\text{Since } E[\mu_{x|y}] = E X = \mu_x$$

13

Exa: CONDITIONAL MEAN & VARIANCE FOR MULTINOMIAL

$$P(X_1 = r_1, X_2 = r_2) = \frac{n!}{r_1! r_2! (n-r_1-r_2)!} p_1^{r_1} p_2^{r_2} (1-p_1-p_2)^{n-r_1-r_2}$$

We saw earlier that

$$P[X_2 = r_2 | X_1 = k] = \frac{(n-k)!}{r_2! (n-k-r_2)!} \left(\frac{p_2}{1-p_1}\right)^{r_2} \left(\frac{1-p_1-p_2}{1-p_1}\right)^{n-k-r_2}$$

$$\sim b(n-k, \frac{p_2}{1-p_1})$$

Thus it is clear that

$$\begin{aligned} E[X_2 | X_1 = k] &= \text{mean of binomial r.v. with} \\ &\quad \text{parameters } (n-k) \text{ and } \frac{p_2}{1-p_1} \\ \text{i.e.} &= (n-k) \frac{p_2}{1-p_1} \\ &= \frac{np_2}{1-p_1} - \left(\frac{p_2}{1-p_1}\right) k \end{aligned}$$

and

$$\begin{aligned} \text{Var}[X_2 | X_1 = k] &= \text{var. of r.var having } b(n-k, \frac{p_2}{1-p_1}) \\ \text{i.e.} &= (n-k) \left(\frac{p_2}{1-p_1}\right) \left(1 - \frac{p_2}{1-p_1}\right) \end{aligned}$$

NOTE : (TERMINOLOGY)

The Conditional expectation  $E(\bar{Y} | x)$  is termed the REGRESSION FUNCTION of  $\bar{Y}$  on  $\bar{X}$

— and the plot of this function against  $x$  is termed a REGRESSION CURVE

When this function is linear, it is referred to as the REGRESSION LINE.

Check : 
$$\begin{aligned} E[E(x_2 | X_1)] &= \frac{np_2}{1-p_1} - \left(\frac{p_2}{1-p_1}\right) E(k) \\ &= \left(\frac{np_2}{1-p_1}\right) - \left(\frac{p_2}{1-p_1}\right) np_1 = \frac{np_2}{1-p_1} [1-p_1] \\ &= np_2 \end{aligned}$$