

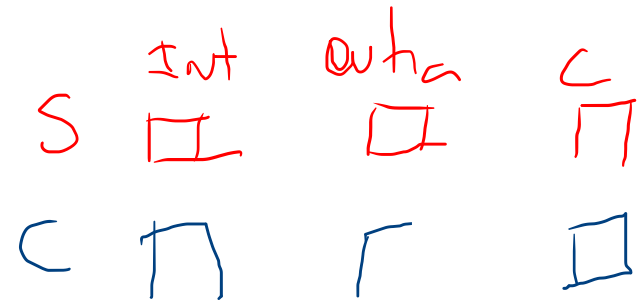
1) EXPLORAÇÃO DOS DADOS NO DB Cap. 4

DADOS JUNTOS AO RH DA EMPRESA //

Nº	Estado civil X_1	Grau de instrução X_2	Nº de filhos X_3	Salário (x sal. mín) X_4	Idade X_5		Região de procedência X_6
					anos	meses	
1	solteiro	ensino fundamental	—	4,00	26	03	interior
2	casado	ensino fundamental	1	4,56	32	10	capital
3	casado	ensino fundamental	2	5,25	36	05	capital
4	solteiro	ensino médio	—	5,73	20	10	outra
5	solteiro	ensino fundamental	—	6,26	40	07	outra
6	casado	ensino fundamental	0	6,66	28	00	interior
7	solteiro	ensino fundamental	—	6,86	41	00	interior
8	solteiro	ensino fundamental	—	7,39	43	04	capital
9	casado	ensino médio	1	7,59	34	10	capital
10	solteiro	ensino médio	—	7,44	23	06	outra
11	casado	ensino médio	2	8,12	33	06	interior
12	solteiro	ensino fundamental	—	8,46	27	11	capital
13	solteiro	ensino médio	—	8,74	37	05	outra
14	casado	ensino fundamental	3	8,95	44	02	outra
15	casado	ensino médio	0	9,13	30	05	interior
16	solteiro	ensino médio	—	9,35	38	08	outra
17	casado	ensino médio	1	9,77	31	07	capital
18	casado	ensino fundamental	2	9,80	39	07	outra
19	solteiro	superior	—	10,53	25	08	interior
20	solteiro	ensino médio	—	10,76	37	04	interior

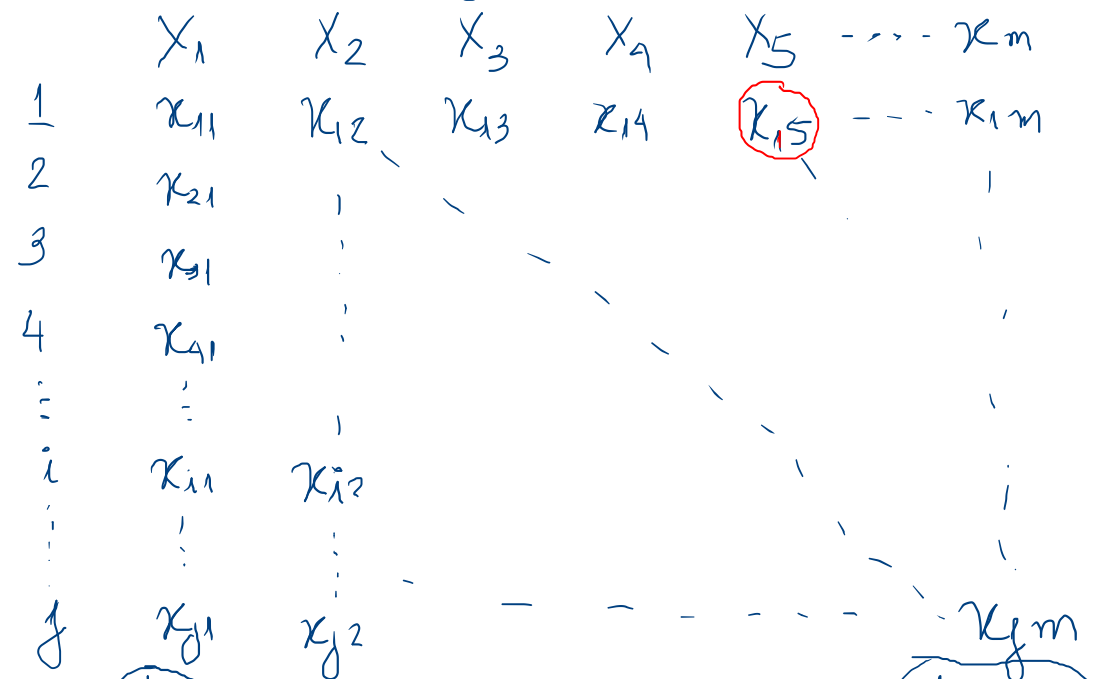
RELEMBRANDO: Variável

- QUALITATIVA
 - ORDINAL: Escola, Idade
 - NOMINAL: Estado civil, Região
- QUANTITATIVA: nº filhos, Idade, salário



Posso fazer:

- $X_1 = \text{Est. civil}$
 - $X_2 = \text{Escolaridade}$
 - $X_3 = \text{Nº filhos}$
 - $X_4 = \text{Idade (anos)}$
 - $X_5 = \text{Procedência}$
 - $X_6 = \text{salário}$
- } → genérico



$$\sum_{t=1}^j x_{tL}$$

$$\sum_{t=1}^j \sum_{k=1}^m x_{tk}$$

↳ somatória facilitar entender nossas contas

2) P/ ANALISAR O BD \rightarrow precisamos investigar

A RELAÇÃO ENTRE NOSSAS VARIÁVEIS.

\Rightarrow Aqui temos que considerar, enfim se as variáveis são DEPENDENTES ou INDEPENDENTES e se são Qualitativas (ordinal/nominal) ou Quantitativa (discretas/contínuas)

Situação 01: (Qualitativa e Qualitativa)

Exemplo: Região de Proveniência \times Estado civil

	X_5			X_1	
$X_1 \backslash X_5$	capital	Interior	Outra	TOTAL	
Solteiro	3	4	4	11	
Casado	4	3	2	9	
	7	7	6	20	

Distribuição conjunta entre X_1 e X_5
dist(x,y)

Posso montar tabelas \neq s c/ os dados

\rightarrow Proporção em relação ao total geral

\rightarrow Proporção a uma variável específica
 \hookrightarrow Linhas \checkmark
 coluna \checkmark

	C	I	O	T	
S	3/20 (15%)	4/20 (20%)	4/20 (20%)	11/20	} Em relação ao total geral
C	4/20	3/20	2/20	9/20	
	7/20	7/20	6/20	20/20	
	C	I	O		
S	3/7	4/7	4/6	11/20	
C	4/7	3/7	2/6	9/20	
	100%	100%	100%	100%	

As variáveis são dependentes ou independentes?

3) Vamos pensar num caso prático

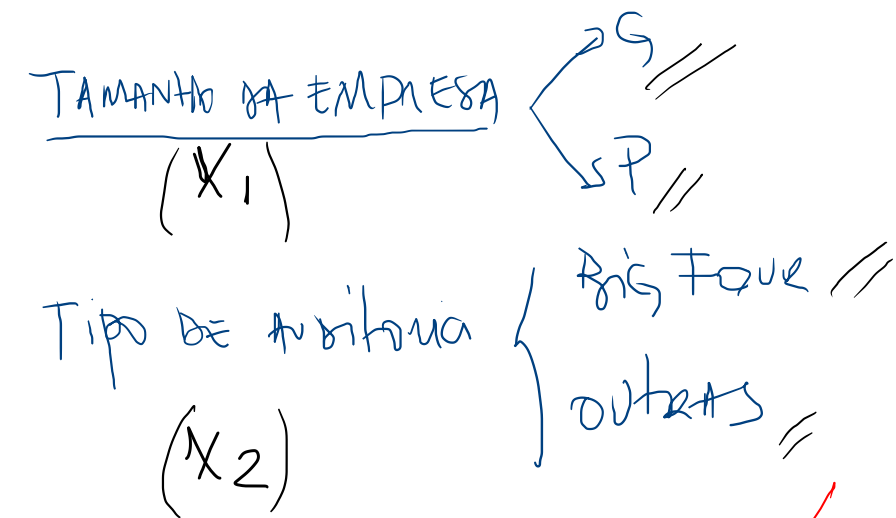


TABELA RESUMO

$X_1 \backslash X_2$	Big Four	outras	TOTAL
G	100 (71%)	20 (33%)	120 (60%)
P	40 (29%)	40 (67%)	80 (40%)
	140 (100%)	60 (100%)	200 (100%)

PARECE HAVER ASSOCIAÇÃO //

BAIXA FREQ. NAS PEQUENAS
ALTA FREQ. NAS GRANDES

COMO MEDIR ESSA ASSOCIAÇÃO?

$C \equiv$ coeficiente de contingência (de Pearson)
or
 $\tilde{C} \equiv$ coeficiente de contingência modificado

Valores nos dados

observado

$X_1 \backslash X_2$	Big Four	Outra	total
G	<u>100</u> (83%)	20 (17%)	120 (100%)
P	40 (50%)	40 (50%)	80 (100%)
TOTAL	140 (<u>70%</u>)	60 (30%)	200 (100%)

LinHA

A) Se n houvesse dependência, empresas
 cerca de 70% de G e P contratariam
 Big farm e 30% outras. Isto é:

→ Esperado

	Big	Outros	TOTAL
G	84 (70%)	36 (30%)	120 100%
P	56 (70%)	24 (30%)	80 100%
TOTAL	140 (70%)	60 (30%)	200 100%

Desvio entre obs e esperado

	Big	OUTRA	$e_i = (o_i - e_e)$
G	$100 - 84 = 16 (3,05)$	$20 - 36 = -16 (7,11)$	
P	$40 - 56 = -16$	$40 - 24 = 16 (10,6)$	

Veja que $\sum_{i=1}^4 e_i = 0 \rightarrow (4,57)$

$$\frac{(O_i - O_e)^2}{O_e} = \frac{(16)^2}{84} = 3,05$$

$\frac{(100 - 84)^2}{84} = 3,05$

Fazendo: $3,05 + 7,11 + 4,57 + 10,6 = 25,33$

χ^2 (qui-quadrado) de Pearson

→ Valores gdes de $\chi^2 \Rightarrow$ Associação

Agencia gerencia C e \tilde{C}

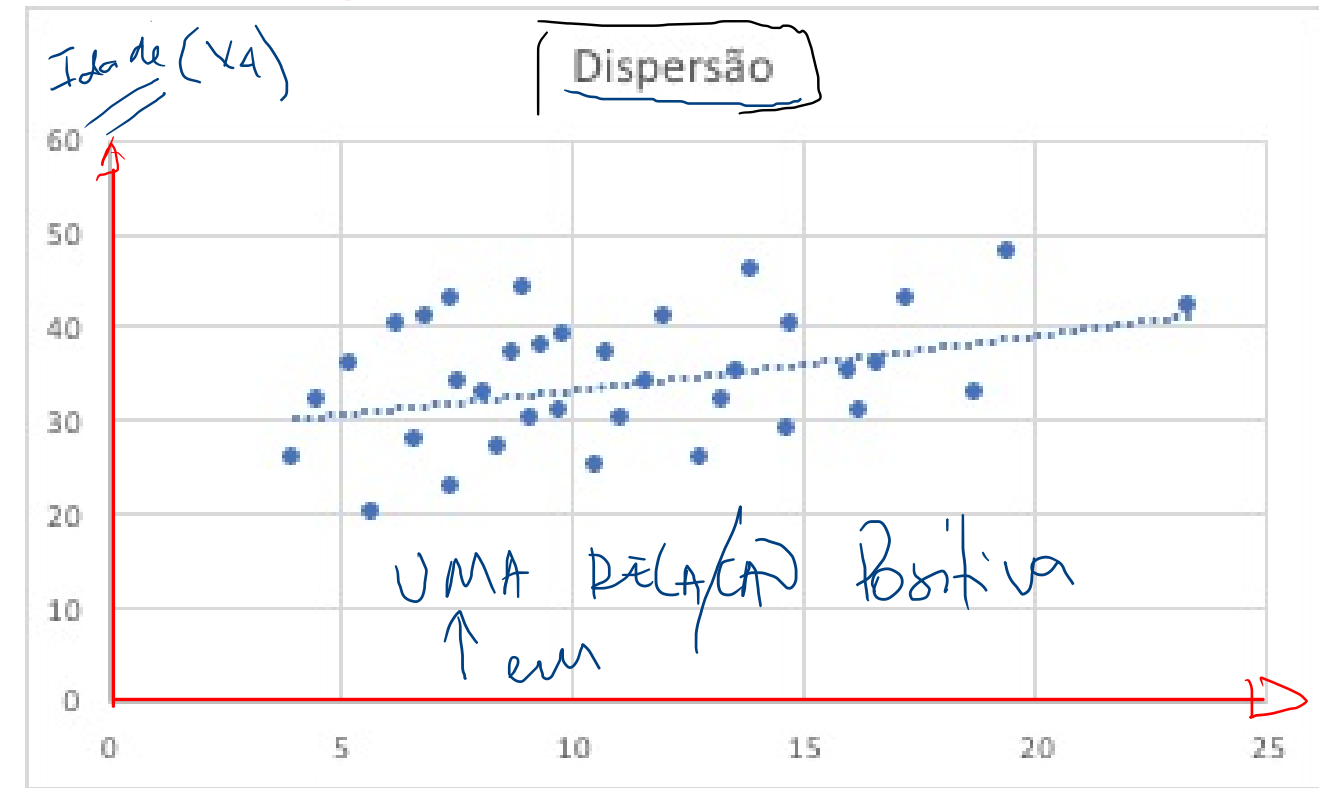
$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}} \approx 0,33 \quad \left| \sqrt{\frac{25,33}{25,33 + 200}}\right.$$

$$\tilde{C} = \sqrt{\frac{\chi^2/n}{(n-1)(s-1)}} \approx 0,36$$

$n = n^\circ$ de linhas (2)
 $s = n^\circ$ de colunas (2)

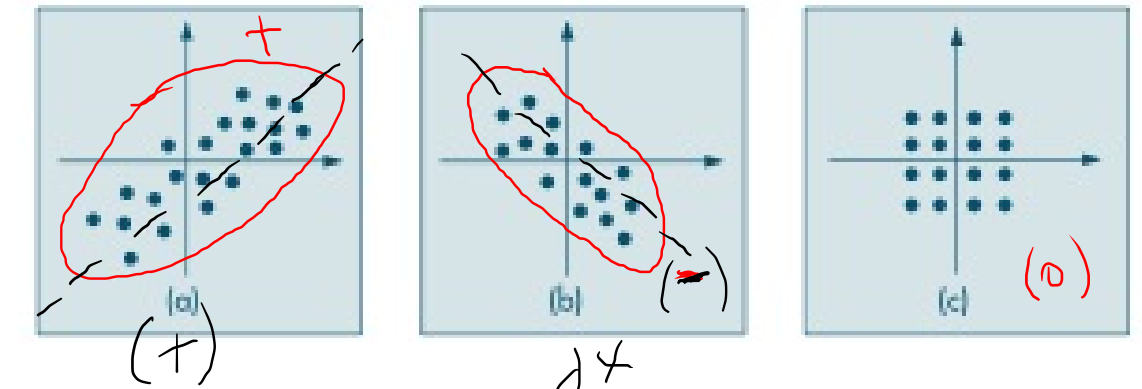
5) Se tivermos X_1 e X_2 quantitativas?

Podemos fazer uma relação entre Idade e salário.



↳ usamos a planilha q. vs usaram p/ fazer a lista //

RELAÇÕES POSSÍVEIS $\frac{dx}{dy} > 0$



A medida tradicional é o $\frac{dx}{dy}$ ind. de correlação

$$r_{(x,y)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{dp(x)} \right) \left(\frac{y_i - \bar{y}}{dp(y)} \right)$$

c/ $(-1 < r < 1)$

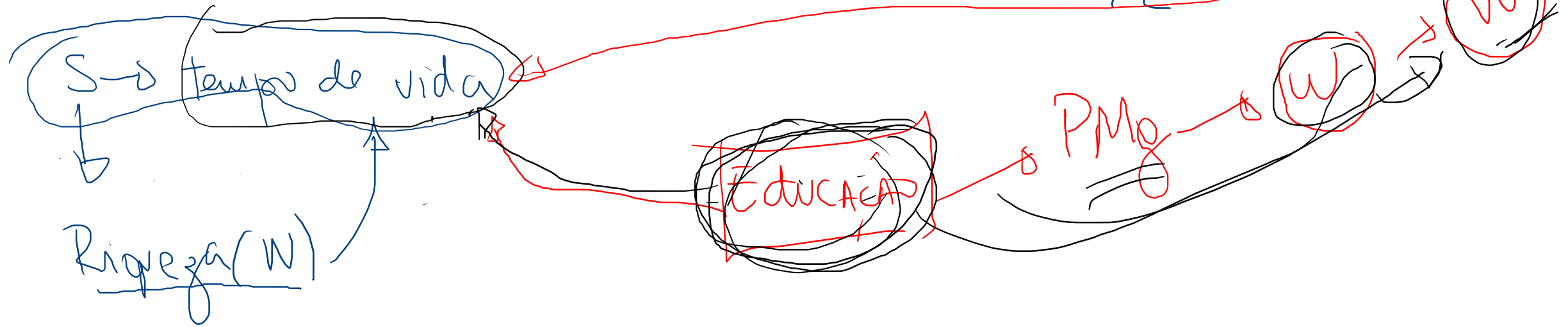
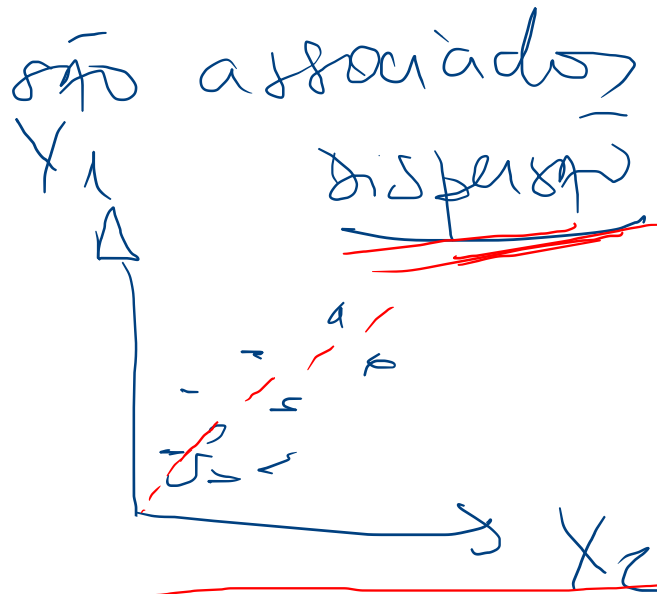
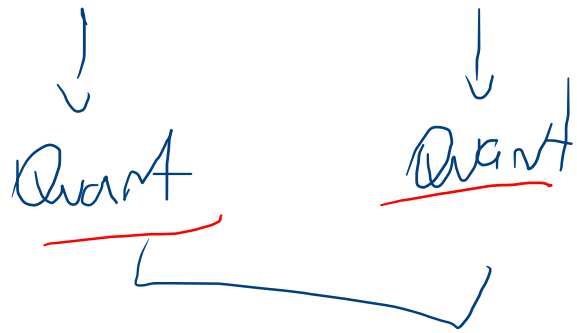
Veja que precisamos dos seguintes dados

\bar{x}_4 ; \bar{x}_6 ; $dp(x_4)$; $dp(x_6)$; n ; $(x_{4i} - \bar{x}_4)$; $(x_{6i} - \bar{x}_6)$

p/ os dados da planilha = $r \approx 0,3$

Raswala

Se Salário e Idade são associados



b) P/ os dados abaixo, podemos fazer:

Vendedor	tempo de casa (x)	Vendas (\$) (y)
1	2	48
2	3	50
3	4	56
4	5	52
5	4	43
6	6	60
7	7	62
8	8	58
9	8	64
10	10	72

$$\sum_{i=1}^{10} x_i = 57$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{10} x_i = 5,7$$

$$\sum_{i=1}^{10} y_i = 565$$

$$\bar{y} = 56,5$$

$\rho = 0,87$

$\frac{x_i - \bar{x}}{\sigma_P(x)}$ ← z_x $\frac{y_i - \bar{y}}{\sigma_P(y)}$ z_y

$(x - \bar{x})$	$(y - \bar{y})$	$\frac{x_i - \bar{x}}{\sigma_P(x)}$	$\frac{y_i - \bar{y}}{\sigma_P(y)}$	$z_x \cdot z_y$
(2-5,7) = -3,7	-0,5	$\frac{-3,7}{2,41} = -1,54$	-1,05	1,617
= -2,7	-6,5	-1,12	-0,8	0,846
= -1,7	-0,5	-0,71	-0,06	0,043
= -0,7	-4,5	-0,29	-0,55	0,160
= -1,7	-13,5	-0,71	-1,66	1,179
= 0,3	3,5	0,12	0,43	0,052
= 1,3	5,5	0,54	0,68	0,367
= 2,3	1,5	0,95	0,19	0,181
= 2,3	7,5	0,95	0,92	0,874
= 4,3	15,5	1,78	1,91	3,4

$Var(x) = \frac{\sum_{i=1}^{10} (x_i - \bar{x})^2}{n} \Rightarrow \sigma_P(x) = 2,41$
 $\sigma_P(y) = 8,11$

$\sum_{i=1}^{10} z_x \cdot z_y = 8,769/10$

7) Outra forma alternativa:

$$\rho(x, y) = \frac{1}{n} \sum \left(\frac{x_i - \bar{x}}{dP(x)} \right) \left(\frac{y_i - \bar{y}}{dP(y)} \right)$$

$$= \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{\left(\sum_1^n x_i^2 - n \bar{x}^2 \right) \left(\sum_1^n y_i^2 - n \bar{y}^2 \right)}}$$

Outra medida: $\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}$
covariância

Assim:

$$\rho = \frac{\text{cov}(x, y)}{dP(x) \cdot dP(y)}$$

