

LIENS ENTRE UNE VARIABLE
QUALITATIVE (e.g. REGION) ET UNE
VARIABLE QUANTITATIVE (e.g. PRIX
D'UN PRODUIT)

P: le prix (en €) est une variable C^0 .

R: la région a 3 modalités: IdF, SO et L-R
↑ ↑ ↑
 Id.-de-France Sud-Ouest Pays-de-l'

Tableau de données

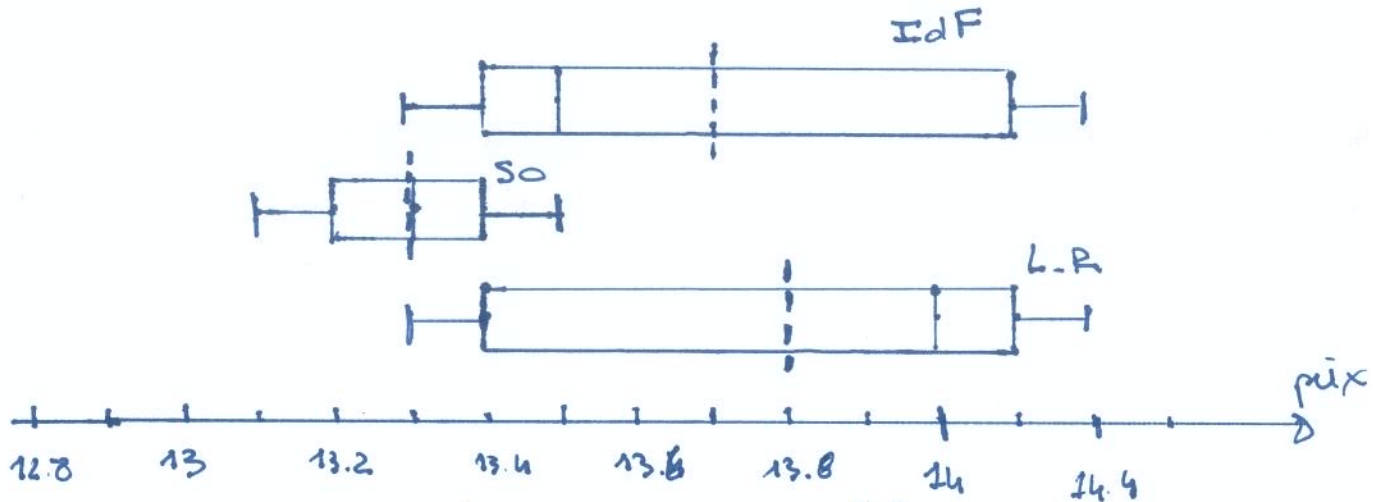
observations	P	R
1	13.5	IdF
2	13.2	SO
3	13.4	L-R
4	14.2	IdF
5	13.3	SO
6	13.3	L-R
7	14.1	IdF
8	13.1	SO
9	14.0	L-R
10	13.4	IdF
11	13.5	SO
12	14.2	L-R
13	13.3	IdF
14	13.4	SO
15	14.1	L-R

↙ Distributions
conditionnelles à
la région

Prix	Région		
	IdF	SO	L-R
Moyen	13.7	13.3	13.8
$Q_{25\%}$	13.4	13.2	13.4
Median	13.5	13.3	14
$Q_{75\%}$	14.1	13.4	14.1
Variances	0.14	0.02	0.14

(.../...)

On peut représenter ces distributions conditionnelles:



⇒ Il y a bien une dispersion prix / région.

On va décomposer la variance totale

$$s_p^2 = \frac{1}{n} \sum (p_i - \bar{p})^2 \quad \text{ou} \quad \bar{p} = \frac{1}{n} \sum_{i=1}^n p_i$$

$$= \frac{1}{n} \sum_{j \in \text{IdF}} p_j + \frac{1}{n_{S.O.}} \sum_{j \in S.O.} p_j + \dots$$

en une variance intra-région et une variance

intra-région

$$V_{\text{intra}} = \frac{1}{n} \sum_{j \in \{\text{IdF}, S.O., L.R.\}} n_{ij} s_{ij}^2 \quad \text{ou} \quad s_{ij}^2 = \frac{1}{n_{ij}} \sum_{i \in \text{région } j} (p_i - p_{ij})^2$$

est la variance de la région "j" et n_{ij} est l'effectif des observations de la région "j"

$$\text{et } V_{\text{intra}} = \frac{1}{n} \sum_{j \in \{\text{IdF}, S.O., L.R.\}} n_{ij} (\bar{p}_{ij} - \bar{p})^2$$

Dans notre exemple numérique : $V_{\text{intra}} = 0.04667$

$$V_{\text{intra}} = 0.1$$

$$s_p^2 = 0.14667$$

Le rapport de corrélation est $\eta_{\text{PIR}} := \sqrt{\frac{V_{\text{intra}}}{s_p^2}} = \sqrt{1 - \frac{V_{\text{intra}}}{s_p^2}}$

Plus ce rapport est proche de 1, plus (...)

les groupes sont hétérogènes et donc il y a une liaison forte entre la variable qualitative (Région) et quantitative (Prix). Si μ_{PIR} est proche de 0, les moyennes conditionnelles fluctent peu par rapport à la moyenne globale : peu de lien entre P et R.

→ Dans notre application numérique, $\mu_{PIR} = 0.564$

permet de conclure qu'il y a une différence selon les régions; l'examen des graphiques "boîtes à moustache" montre bien que c'est la région SO qui est franchement différente des 2 autres.

On verra au passage l'asymétrie (moyenne très différente de la médiane) des distributions conditionnelles des prix en

IdF et L-R

⑨ Démonstration de $S_p^2 = V_{inter} + V_{intra}$:

$$S_p^2 = \frac{1}{n} \sum_{i=1}^n (P_i - \bar{P})^2 = \frac{1}{n} \left[\sum_{i \in \text{IdF}} (P_i - \bar{P})^2 + \sum_{i \in \text{SO}} (P_i - \bar{P})^2 + \sum_{i \in \text{L-R}} (P_i - \bar{P})^2 \right]$$

= idem

$$= \frac{1}{n} \left[(P_i - \bar{P}_{\text{IdF}})^2 + (\bar{P}_{\text{IdF}} - \bar{P})^2 + 2(P_i - \bar{P}_{\text{IdF}})(\bar{P}_{\text{IdF}} - \bar{P}) \right]$$

= idem

$\forall i, j \in \{\text{IdF}, \text{SO}, \text{L-R}\}$

$$\hookrightarrow \sum_{i \in j} (P_i - \bar{P})^2 = \sum_{i \in j} \left[(P_i - \bar{P}_{j'})^2 + (\bar{P}_{j'} - \bar{P})^2 + 2(P_i - \bar{P}_{j'}) (\bar{P}_{j'} - \bar{P}) \right]$$

$$= n_{j'} \cdot \frac{1}{n_{j'}} \sum_{i \in j} (P_i - \bar{P}_{j'})^2 + n_{j'} (\bar{P}_{j'} - \bar{P})^2 + 0$$

= $S_{j'}^2$

, car $\sum_{i \in j} (P_i - \bar{P}_{j'}) = 0$

En rassemblant les termes : = V_{inter}

= V_{intra}

$$S_p^2 = \frac{1}{n} \sum_{j'} n_{j'} \cdot S_{j'}^2 + \frac{1}{n} \sum_{j'} n_{j'} (\bar{P}_{j'} - \bar{P})^2 \quad \square$$