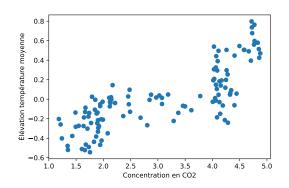
Régression Linéaire Simple

Benoit Gaüzère, Stéphane Canu benoit.gauzere@insa-rouen.fr

INSA Rouen Normandie - ITI

March 20, 2025

Observation de deux variables



Statistiques descriptives

- Analyse univariée
- Analyse bivariée : dépendance linéaire

Comment expliciter cette relation?

Le problème de la régression linéaire

Les données

n couple d'observations:

- $ightharpoonup x_i \in \mathbb{R}$: la variable explicative
- $ightharpoonup y_i \in {\rm I\!R}$: la variable à expliquer, prédire

Le problème

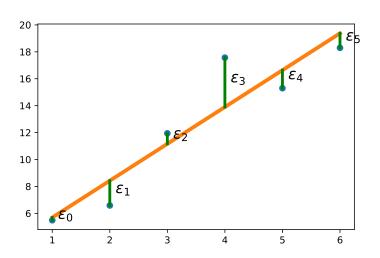
Trouver la droite qui représente au mieux la relation linéaire

$$\widehat{y} = f(x) \simeq ax + b$$

Mais les points ne sont pas forcément sur la droite

$$\widehat{y}_i = ax_i + b + \varepsilon_i$$

Modélise le bruit



- $ightharpoonup a \simeq 2,73$
- $b \simeq 2.95$
- $\epsilon_i \simeq [-0.197, -1.840, 0.781, 3.670, -1.339, -1.076]$

Le Modèle Linéaire

Définition : Modèle Linéaire

Le modèle linéaire pose la relation suivante entre la variable explicative x et la variable à expliquer y avec les paramètres inconnus (a,b,ε)

$$y = ax + b + \varepsilon$$
 avec $\mathbf{a} = (a, b) \in \mathbb{R}^2$

- Hypothèse : Observations = modèle + bruit
- $\triangleright \varepsilon \sim \mathcal{N}(0, \sigma^2)$

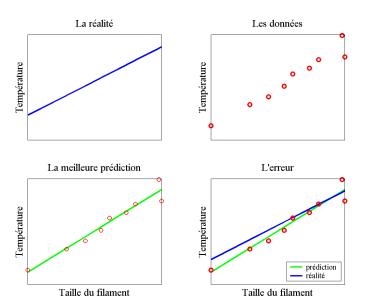
Résumé du vocabulaire

| $lackbox{ Variables explicatives } \ldots x \in {\rm I\!R}$ |
|---------------------------------------------------------------------------------------------------|
| $lackbox{ Variable à expliquer } \dots y \in {\rm I\!R}$ |
| lacktriangle Erreur ou bruit $arepsilon$ |
| $lackbox{ Paramètres scalaires} \ldots a,b \in {\rm I\!R}$, |
| $lackbox{Paramètres}$ (forme vectorielle) $\dots \dots \mathbf{a} \in \mathbb{R}^2$ |
| $\blacktriangleright \ Mod\`{ele} \ \ldots \dots y = f(\mathbf{x}, \mathbf{a}) + \varepsilon$ |
| lacktriangle Estimation |
| \blacktriangleright Prédiction |
| $lackbox{ Variables aléatoires } \ldots \varepsilon$ et donc y et donc a^\star |
| |

$$y = \underbrace{ax + b}_{f(x,\mathbf{a})} + \varepsilon$$

Un exemple : l'étalonnage d'un capteur

Les différentes phases de la régression



Les moindres carrés pour la régression simple

Moindres carrés

Fonction objectif

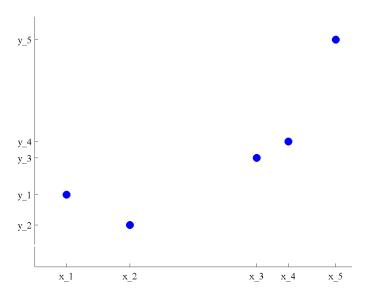
$$\min_{a,b} J(a,b)$$

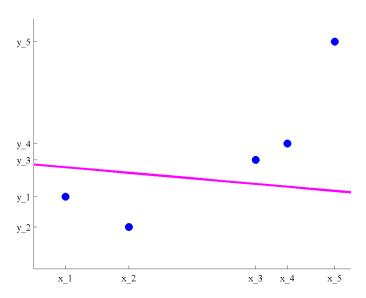
avec

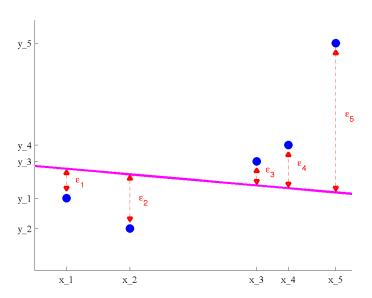
$$J(a,b) = \sum_{i=1}^{n} \left(\underbrace{y_i - ax_i - b}_{\varepsilon_i} \right)^2$$

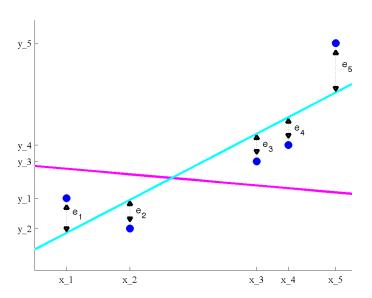
Interprétation

Ce problème peut s'interpréter comme la recherche de la droite d'équation ax+b passant "au mieux" (au sens des moindres carrés) parmi le nuage des observations (x_i,y_i) , $i=1,\ldots,n$.









Le problème des moindres carrés

Les données dont nous disposons peuvent aussi être vue comme un système de n équations à 2+n inconnues $(a, b \text{ et les } \varepsilon_i)$. Ce système s'écrit de la manière suivante :

$$\begin{cases} ax_1 + b + \varepsilon_1 &= y_1 \\ \vdots & \vdots \\ ax_i + b + \varepsilon_i &= y_i \\ \vdots & \vdots \\ ax_n + b + \varepsilon_n &= y_n \end{cases}$$

On recherche a et b qui minimisent simultanément tous les ε_i

$$J(a,b) = \sum_{i=1}^{n} \left(\underbrace{y_i - ax_i - b}_{\varepsilon_i} \right)^2$$

Calcul du gradient

$$\underset{a,b}{\operatorname{arg \, min}} J(a,b)$$
 avec $J(a,b) = \frac{1}{2} \sum_{i=1}^{n} (ax_i + b - y_i)^2$

Méthode du gradient

 \blacktriangleright (a^{\star}, b^{\star}) est solution du problème

$$\underset{a,b}{\operatorname{arg\,min}} J(a,b) \quad \Leftrightarrow \begin{cases} \frac{\partial J(a^{\star},b^{\star})}{\partial a} = 0\\ \frac{\partial J(a^{\star},b^{\star})}{\partial b} = 0 \end{cases}$$

Dérivées partielles de J(a, b)

$$\frac{\partial J(a,b)}{\partial a} = \sum_{i=1}^{n} (ax_i + b - y_i) x_i = \sum_{i=1}^{n} (ax_i^2 + bx_i - y_i x_i)$$

$$= \sum_{i=1}^{n} ax_i^2 + \sum_{i=1}^{n} bx_i - \sum_{i=1}^{n} y_i x_i$$

$$= a \sum_{i=1}^{n} x_i^2 + b \sum_{i=1}^{n} x_i - \sum_{i=1}^{n} y_i x_i$$

$$\frac{\partial J(a,b)}{\partial b} = \sum_{i=1}^{n} (ax_i + b - y_i) = \sum_{i=1}^{n} ax_i + \sum_{i=1}^{n} b - \sum_{i=1}^{n} y_i$$
$$= a \sum_{i=1}^{n} x_i + bn - \sum_{i=1}^{n} y_i$$

Calcul de a^* et b^*

Deux équations linéaires à deux inconnues

$$\begin{cases}
\frac{\partial J(a^{\star}, b^{\star})}{\partial a} = 0 & a^{\star} \sum_{i=1}^{n} x_i^2 + b^{\star} \sum_{i=1}^{n} x_i = \sum_{i=1}^{n} y_i x_i(1) \\
\frac{\partial J(a^{\star}, b^{\star})}{\partial b} = 0 & a^{\star} \sum_{i=1}^{n} x_i + b^{\star} n = \sum_{i=1}^{n} y_i \quad (2)
\end{cases}$$

Calcul de a^{\star} et b^{\star}

Calcul de a^*

$$ightharpoonup (1) * n - (2) * \sum_{i=1}^{n} x_i$$

$$a^* = \frac{n \sum_{i=1}^{n} y_i x_i - \sum_{i=1}^{n} y_i \sum_{i=1}^{n} x_i}{n \sum_{i=1}^{n} x_i^2 - \left(\sum_{i=1}^{n} x_i\right)^2}$$

Calcul_de
$$b^*$$

$$\star = \frac{\sum_{i=1}^{n} y_i - a^* \sum_{i=1}^{n} x_i}{n}$$

Réécrivons la solution

Soit
$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
, $\overline{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$ et $V_x = \frac{1}{n} \sum_{i=1}^{n} (x_i - \overline{x})^2$.

$$a^{\star} = \frac{n \sum_{i=1}^{n} y_{i} x_{i} - \sum_{i=1}^{n} y_{i} \sum_{i=1}^{n} x_{i}}{n \sum_{i=1}^{n} x_{i}^{2} - \left(\sum_{i=1}^{n} x_{i}\right)^{2}} = \frac{\frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - \frac{1}{n} \sum_{i=1}^{n} y_{i} \frac{1}{n} \sum_{i=1}^{n} x_{i}}{\frac{1}{n} \sum_{i=1}^{n} x_{i}^{2} - \left(\frac{1}{n} \sum_{i=1}^{n} x_{i}\right)^{2}}$$

$$= \frac{\frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - \overline{y} \overline{x}}{\frac{1}{n} \sum_{i=1}^{n} x_{i}^{2} - \overline{x}^{2}} = \frac{\frac{1}{n} \sum_{i=1}^{n} (y_{i} - \overline{y}) (x_{i} - \overline{x})}{\frac{1}{n} \sum_{i=1}^{n} (x_{i} - \overline{x})^{2}}$$

$$a^{\star} = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{V}$$

$$b^* = \frac{\sum_{i=1}^n y_i - a^* \sum_{i=1}^n x_i}{n}$$
$$= \frac{1}{n} \sum_{i=1}^n y_i - a^* \frac{1}{n} \sum_{i=1}^n x_i$$
$$= \overline{y} - a^* \overline{x}$$

Le théorème des moindres carrés

Théorème : Théorème des moindres carrées

Soit (x_i, y_i) , i = 1, n un ensemble de couples d'observations.

La solution du problème de minimisation de la somme des carrées des erreurs

$$\min_{a,b} \sum_{i=1}^{n} (ax_i + b - y_i)^2$$

est donnée par a* et b* définis par :

$$a^{\star} = \frac{\sum_{i=1}^{n} (y_i - \overline{y}) (x_i - \overline{x})}{\sum_{i=1}^{n} (x_i - \overline{x})^2} = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{V_x} \text{ et } b^{\star} = \overline{y} - a^{\star} \overline{x}$$

 a^{\star} et b^{\star} sont les estimateurs au sens des moindres carrés.

En Résumé

Apprentissage des paramètres

- ightharpoonup Dépend de x et y
- $\qquad \qquad \bullet \quad a^\star = \frac{\mathsf{cov}(\mathbf{x},\mathbf{y})}{V_r}$
- $b^* = \overline{y} a^* \ \overline{x}$

Modèle prédictif

- ▶ Dépend d'une observation x, et des paramètres a^* et b^*
- \triangleright Prédire une nouvelle observation x:

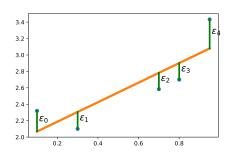
$$\widehat{y} = f(x) = a^{\star}x + b^{\star}$$

• Erreur de prédiction $\varepsilon = y_i - \widehat{y}_i$

Mise en oeuvre

e = y - yp

```
mx = np.mean(x)
my = np.mean(y)
sxx = np.sum((x-mx)**2)
sxy = (x-mx).T @ (y-my)
a = sxy/sxx
b = my - a*mx
yp = a*x+b
```



| x | y | y_p | $arepsilon_i$ |
|------|--------|--------|---------------|
| 0.1 | 2.3196 | 2.0681 | 0.2514 |
| 0.3 | 2.1000 | 2.3061 | -0.2061 |
| 0.7 | 2.5836 | 2.7820 | -0.1984 |
| 0.8 | 2.7000 | 2.9010 | -0.2010 |
| 0.95 | 3.4335 | 3.0795 | 0.3540 |

Le poids des observations

Influence des observations dans le calcul de a^*

$$a^{\star} = \frac{\sum_{i=1}^{n} (x_{i} - \bar{x}) (y_{i} - \bar{y})}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} = \frac{\sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}$$

$$= \sum_{i=1}^{n} \frac{x_{i} - \bar{x}}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} y_{i}$$

$$= \sum_{i=1}^{n} g(x_{i}) y_{i} = \sum_{i=1}^{n} w_{i} y_{i}$$

$$w_i = \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{ mesure le poids de l'observation } x_i \text{ dans }$$
 le calcul de a^\star .

Remarques

$$a^{\star} \text{ et corrélation} \\ a^{\star} = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{s_x^2} = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{s_x s_x} \frac{s_y}{s_y} = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{s_x s_y} \frac{s_y}{s_x} = \text{cor}(\mathbf{x}, \mathbf{y}) \frac{s_y}{s_x}$$

Droite de régression

$$\begin{array}{rcl} y & = & a^{\star}x + b^{\star} \\ & = & a^{\star}x + \bar{y} - a^{\star} \; \bar{x} \\ & = & a^{\star}(x - \bar{x}) + \bar{y} & = \frac{\mathsf{COV}(\mathbf{x}, \mathbf{y})}{s_x^2}(x - \bar{x}) + \bar{y} \end{array}$$

La droite de régression passe par le point (\bar{x}, \bar{y}) :

$$f(\bar{x}) = a^* \bar{x} + b^* = \frac{\operatorname{cov}(\mathbf{x}, \mathbf{y})}{s_x^2} \underbrace{(x - \bar{x})}_{=0} + \bar{y}$$

Coefficient de détermination R^2

Quelle quantité de la relation est expliquée ?

- Écart expliqué par le modèle : $\sum_{i=1}^{n} (\widehat{y}_i \overline{y})^2$
- Écart total : $\sum_{i=1}^{n} (y_i \bar{y})^2$
- Écart résiduel : $\sum_{i=1}^{n} (y_i \widehat{y}_i)^2$

 \mathbb{R}^2

Le coefficient de détermination $R^2 \in \{0,1\}$ est le rapport de l'écart expliqué versus l'écart total

$$R^{2} = \frac{\sum_{i=1}^{n} (\widehat{y}_{i} - \bar{y})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$

Prédiction et erreur de prédiction

Erreurs et variable aléatoire

- Les ε_i suivent une loi normale $\mathcal{N}(0, \sigma^2)$
- $\widehat{\sigma^2} = \frac{1}{n-2} \sum_{i=1}^n \varepsilon_i^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i a^* x_i b^*)^2$

y_x est également une v.a.

$$y_x = a^*x + b^* + \varepsilon = \bar{y} + a^*(\mathbf{x} - \bar{x}) + \varepsilon$$

$$\widehat{\mathbf{E}(y_x)} = a^*x + b^* = \bar{y} + a^*(x - \bar{x})$$

$$V(y_x) = \sigma^2 \left(1 + \frac{1}{n} + \frac{(x-\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

Intervalle de prédiction

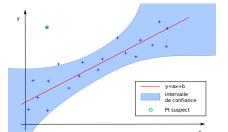
Théorème : Intervalle de prédicition

Pour un x donné, il est probable, avec une probabilité $1-\alpha$, que

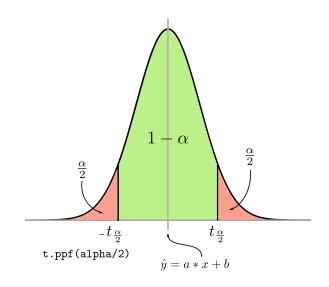
$$y_x \in \left\{ \bar{y} + a^*(x - \bar{x}) \pm t_{\frac{\alpha}{2}} \widehat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right\}$$

où t_{α} vérifie, relativement à la loi de Student à n-2 degrès de libertés, $P(T \leq t_{\alpha}) = \alpha$.

Plus on s'éloigne de la moyenne des x, plus l'incertitude grandit



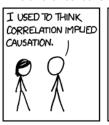
Calculer $t_{\frac{\alpha}{2}}$ avec scipy



Conclusion

La régression simple

- observation = modèle + bruit
- modèlisation linéaire d'une relation
- certains modèles non linéaires peuvent être linéarisés (projection)
- Moindres carrés : facile à calculer







La suite

- ▶ Passer au multi dimensionnel (plusieurs variables explicatives)
- ► Valider le modèle : diagnostic de la régression

Éléments de démonstration de $Var(y_x)$

$$V(a^{\star}) = V\left(\frac{\sum_{i=1}^{n} (y_{i} - \bar{y}) (x_{i} - \bar{x})}{\sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{x})^{2}}\right)$$

$$= \left(\frac{1}{\sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{x})^{2}}\right)^{2} V\left(\sum_{i=1}^{n} (a(x_{i} - \bar{x}) + \varepsilon) (x_{i} - \bar{x})\right)$$

$$= \left(\frac{1}{\sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{x})^{2}}\right)^{2} V\left(\sum_{i=1}^{n} \varepsilon (x_{i} - \bar{x})\right)$$

$$= \frac{\sigma^{2}}{\sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{x})^{2}}$$

D'où
$$V(y_x) = V(\bar{y}) + V(a^\star)(\mathbf{x} - \bar{x})^2 + \sigma^2$$

$$= \frac{\sigma^2}{n} + \frac{\sigma^2}{\sum_{i=1}^n (\mathbf{x}_i - \bar{x})^2} (\mathbf{x} - \bar{x})^2 + \sigma^2$$

$$= \sigma^2 \left(1 + \frac{1}{n} + \frac{(\mathbf{x} - \bar{x})^2}{\sum_{i=1}^n (\mathbf{x}_i - \bar{x})^2} \right)$$

$$= \sigma^2 \left(1 + \frac{1}{n} + \frac{(\mathbf{x} - \bar{x})^2}{\sum_{i=1}^n (\mathbf{x}_i - \bar{x})^2} \right)$$

Dérivée partielle J(a,b) de par rapport à a

$$J(a,b) = \frac{1}{2} \sum_{i=1}^{n} (ax_i + b - y_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^{n} a^2 x_i^2 + bax_i - ax_i y_i + bax_i + b^2 - by_i - y_i ax_i - bx_i$$

$$\frac{\partial J(a,b)}{\partial a} = \sum_{i=1}^{n} ax_i^2 + bx_i - y_i ax_i$$

$$= \sum_{i=1}^{n} (ax_i + b - y_i) x_i$$

Calcul de a^*

$$\begin{array}{l} \bullet \quad (1)*n-(2)*\sum_{i=1}^{n}x_{i} \text{ (pour supprimer } b^{\star}) \\ n(a^{\star}\sum_{i=1}^{n}x_{i}^{2}+b^{\star}\sum_{i=1}^{n}x_{i})-(\sum_{i=1}^{n}x_{i})(a^{\star}\sum_{i=1}^{n}x_{i}+b^{\star}n) = \\ n\sum_{i=1}^{n}y_{i}x_{i}-\sum_{i=1}^{n}x_{i}\sum_{i=1}^{n}y_{i} \\ \Rightarrow a^{\star}n\sum x_{i}^{2}+b^{\star}n\sum x_{i}-a^{\star}(\sum x_{i})^{2}-b^{\star}n\sum x_{i} \\ \Rightarrow a^{\star}(n(\sum x_{i}^{2})-(\sum x_{i})^{2}=n\sum y_{i}x_{i}-\sum x_{i}\sum y_{i} \\ \Rightarrow a^{\star}=\frac{n\sum_{i=1}^{n}y_{i}x_{i}-\sum_{i=1}^{n}y_{i}\sum_{i=1}^{n}x_{i}}{n\sum_{i=1}^{n}x_{i}^{2}-\left(\sum_{i=1}^{n}x_{i}\right)^{2}} \end{array}$$

Calcul de b^*

Connaissant a^* , on reprend (2)

Réécriture de a*

$$\begin{split} a^{\star} &= \frac{n \sum_{i=1}^{n} y_{i} x_{i} - \sum_{i=1}^{n} y_{i} \sum_{i=1}^{n} x_{i}}{n \sum_{i=1}^{n} x_{i}^{2} - \left(\sum_{i=1}^{n} x_{i}\right)^{2}} \Rightarrow * \frac{1}{n^{2}} \\ &= \frac{\frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - \frac{1}{n^{2}} \sum_{i=1}^{n} y_{i} \sum_{i=1}^{n} x_{i}}{\frac{1}{n} \sum_{i=1}^{n} x_{i}^{2} - \frac{1}{n^{2}} \left(\sum_{i=1}^{n} x_{i}\right)^{2}} \\ &= \frac{\frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - \frac{1}{n} \sum_{i=1}^{n} y_{i} \frac{1}{n} \sum_{i=1}^{n} x_{i}}{\frac{1}{n} \sum_{i=1}^{n} x_{i}^{2} - \left(\frac{1}{n} \sum_{i=1}^{n} x_{i}\right)^{2}} \\ &= \frac{\frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - \bar{y}\bar{x}}{\frac{1}{n} \sum_{i=1}^{n} x_{i}^{2} - \bar{x}^{2}} \\ &= \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{V_{x}} \end{split}$$

Théorème de Konig-Huygens

$$(\frac{1}{n}\sum x_i^2) - \bar{x}^2 = \frac{1}{n}\sum (x_i - \bar{x})^2$$

https://fr.wikipedia.org/wiki/Th%C3%A9or%C3%A8me_de_

K%C3%B6nig-Huygens

$$\frac{1}{n}\sum y_i x_i - \bar{y}\bar{x} = \frac{1}{n}\sum (y_i - \bar{y})(x_i - \bar{x})$$

$$\frac{1}{n}\sum_{i}(y_i-\bar{y})(x_i-\bar{x}) = \frac{1}{n}\sum_{i}y_ix_i - y_i\bar{x} - \bar{y}x_i + \bar{y}\bar{x}$$

 $\sum (x_i - \bar{x})(y_i - \bar{y}) = \sum (x_i - \bar{x})y_i$

$$= \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i \bar{x} - \frac{1}{n} \sum \bar{y} x_i + \bar{y} \bar{x}$$

$$= \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i + \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i + \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i + \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i + \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i + \frac{1}{n} \sum y_i x_i - \frac{1}{n} \sum y_i x_i + \frac{1}{n} \sum y_i x_i - \frac{$$

$$= \frac{1}{n} \sum y_i x_i - \bar{x} \frac{1}{n} \sum y_i - \bar{y} \frac{1}{n} \sum x_i + \bar{y} \bar{x}$$

$$= \frac{1}{n} \sum y_i x_i - \bar{x} \bar{y} - \bar{x} \bar{y} + \bar{y} \bar{x}$$

$$egin{align} &=rac{1}{n}\sum y_ix_i-ar{x}rac{1}{n}\sum y_i-ar{y}rac{1}{n}\sum x_i+ar{y} \ &=rac{1}{n}\sum y_ix_i-ar{x}ar{y}-ar{x}ar{y}+ar{y}ar{x} \ \end{pmatrix}$$

 $= \frac{1}{2} \sum y_i x_i - \bar{x}\bar{y}$

$$= \sum x_i y_i - \sum x_i \bar{y} - \sum \bar{x} y_i + \sum \bar{x} \bar{y}$$

$$= \sum x_i y_i - \sum x_i \frac{1}{n} \sum y_j - \sum \frac{1}{n} \sum x_j y_i + \sum \bar{x} \bar{y}$$

$$= \sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i - \frac{1}{n} \sum x_i y_i + \sum \bar{x} \bar{y}$$

$$= \sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i - \frac{1}{n} \sum x_i y_i + \sum \frac{1}{n} \sum x_i y_i - \sum x_i$$

 $\sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y}$

 $=\sum x_iy_i - \bar{x}y_i$

 $=\sum (x_i - \bar{x})y_i$