

Uncorrelation and Independence

Probability and Statistics for Data Science

Carlos Fernandez-Granda



These slides are based on the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

Goal

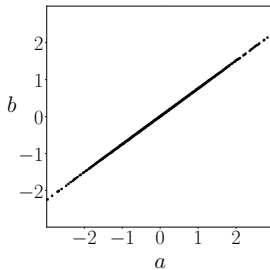
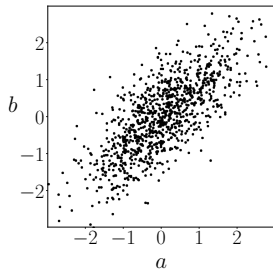
Explain the difference between uncorrelation and independence

$$\rho_{\tilde{a}, \tilde{b}} = 0.75$$

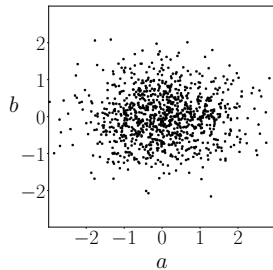
If $\rho_{\tilde{a}, \tilde{b}}$ and $\text{Cov}[\tilde{a}, \tilde{b}]$ are positive, \tilde{a} and \tilde{b} are positively correlated

Linear estimate

$$b = \rho_{\tilde{a}, \tilde{b}} a$$



Residual

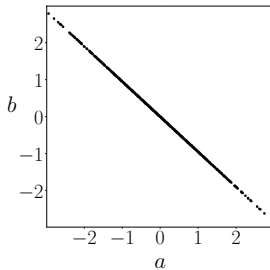
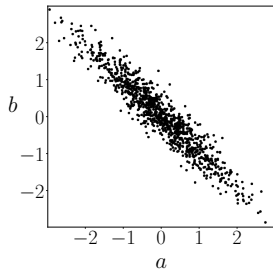


$$\rho_{\tilde{a}, \tilde{b}} = -0.95$$

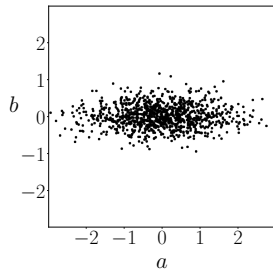
If $\rho_{\tilde{a}, \tilde{b}}$ and $\text{Cov}[\tilde{a}, \tilde{b}]$ are negative, \tilde{a} and \tilde{b} are negatively correlated

Linear estimate

$$\hat{b} = \rho_{\tilde{a}, \tilde{b}} \tilde{a}$$



Residual

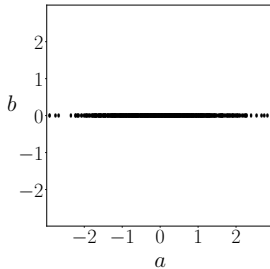
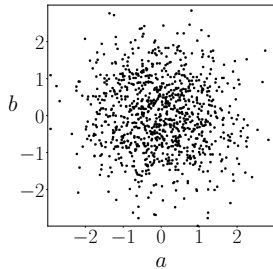


$$\rho_{\tilde{a}, \tilde{b}} = 0$$

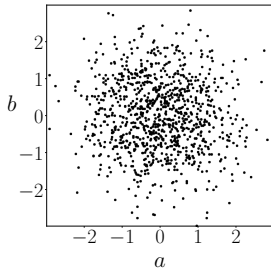
If $\rho_{\tilde{a}, \tilde{b}}$ and $\text{Cov}[\tilde{a}, \tilde{b}]$ are **zero**, \tilde{a} and \tilde{b} are **uncorrelated**

Linear estimate

$$b = \rho_{\tilde{a}, \tilde{b}} a$$



Residual



Independence implies uncorrelation

If \tilde{a} and \tilde{b} are independent, then

$$\begin{aligned}\text{Cov}[\tilde{a}, \tilde{b}] &= \text{E}[\tilde{a}\tilde{b}] - \text{E}[\tilde{a}] \text{E}[\tilde{b}] \\ &= \text{E}[\tilde{a}] \text{E}[\tilde{b}] - \text{E}[\tilde{a}] \text{E}[\tilde{b}] \\ &= 0\end{aligned}$$

Gaussian random variables

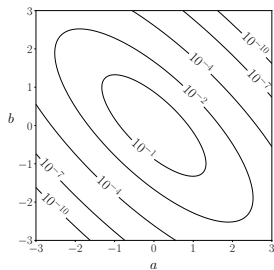
Covariance matrix for uncorrelated Gaussian random variables with zero mean and unit variance

$$\Sigma := \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \Sigma^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

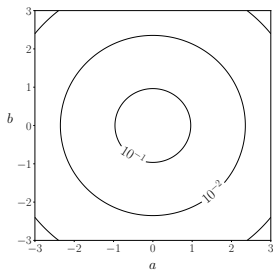
$$\begin{aligned} f_{\tilde{a}, \tilde{b}}(a, b) &= \frac{1}{2\pi \sqrt{|\Sigma|}} \exp \left(-\frac{1}{2} \begin{bmatrix} a \\ b \end{bmatrix}^T \Sigma^{-1} \begin{bmatrix} a \\ b \end{bmatrix} \right) \\ &= \frac{1}{2\pi} \exp \left(-\frac{a^2 + b^2}{2} \right) \\ &= \frac{1}{2\pi} \exp \left(-\frac{a^2}{2} \right) \frac{1}{2\pi} \exp \left(-\frac{b^2}{2} \right) \\ &= f_{\tilde{a}}(a) f_{\tilde{b}}(b) \quad \text{Independent} \end{aligned}$$

Gaussian random variables

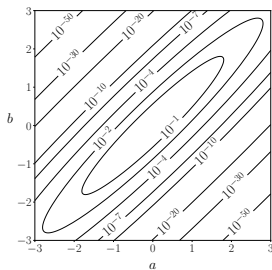
$$\rho = -0.75$$



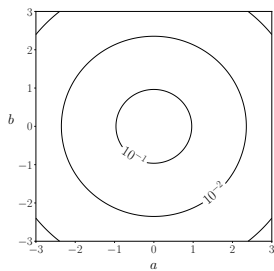
$$\rho = 0$$



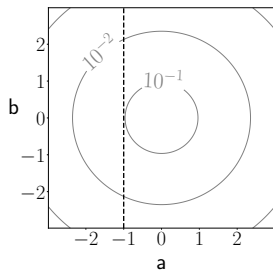
$$\rho = 0.95$$



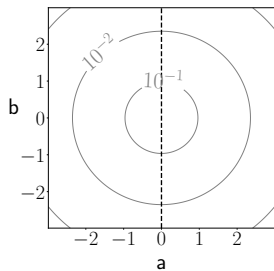
$$\rho = 0$$



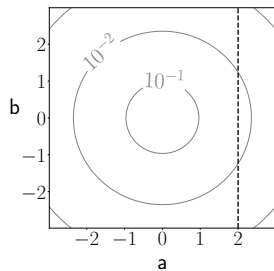
$$a = -1$$



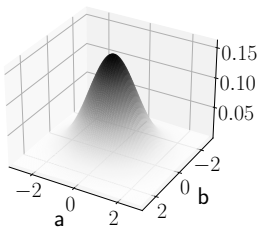
$$a = 0$$



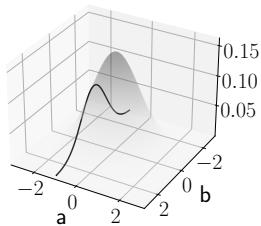
$$a = 2$$



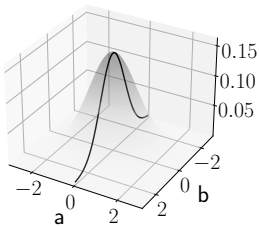
Conditional distribution given $\tilde{a} = a$



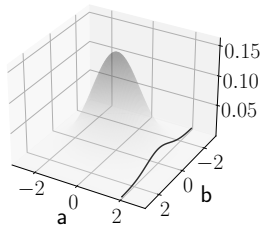
$a = -1$



$a = 0$

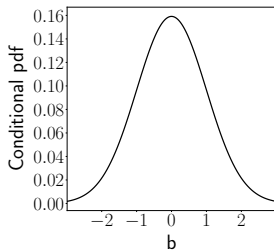


$a = 2$

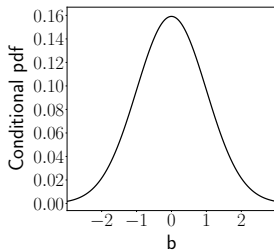


Conditional distribution given $\tilde{a} = a$

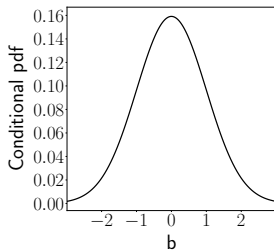
$$a = -1$$



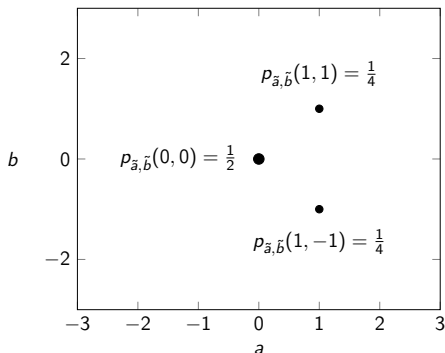
$$a = 0$$



$$a = 2$$



Example

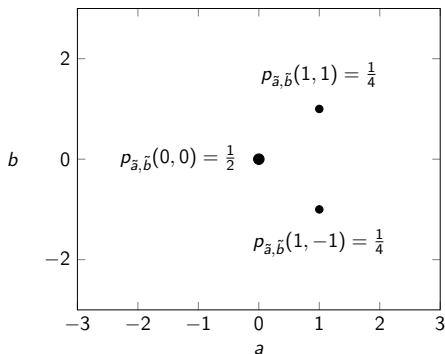


$$E[\tilde{b}] = \sum_{a=0}^1 \sum_{b=-1}^1 b p_{\tilde{a}, \tilde{b}}(a, b) = 0 \cdot \frac{1}{2} - 1 \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} = 0$$

$$\text{Cov}[\tilde{a}, \tilde{b}] = E[\tilde{a}\tilde{b}] - E[\tilde{a}]E[\tilde{b}] = E[\tilde{a}\tilde{b}] \quad \text{Uncorrelated}$$

$$= \sum_{a=0}^1 \sum_{b=-1}^1 ab p_{\tilde{a}, \tilde{b}}(a, b) = 0 \cdot \frac{1}{2} - 1 \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} = 0$$

Example



Conditional pmf of \tilde{b} given $\tilde{a} = 0$?

$$p_{\tilde{b}|\tilde{a}}(0|0) = 1$$

Conditional pmf of \tilde{b} given $\tilde{a} = 1$?

$$p_{\tilde{b}|\tilde{a}}(1|1) = \frac{1}{2} \quad p_{\tilde{b}|\tilde{a}}(-1|1) = \frac{1}{2} \quad \text{Not independent}$$

Uncorrelated residual

Let $\ell_{\text{MMSE}}(\tilde{a})$ be the linear MMSE estimate of \tilde{b} given \tilde{a}

$$\text{Cov} [\tilde{a}, \tilde{b} - \ell_{\text{MMSE}}(\tilde{a})] = 0$$

Height of NBA players

Data:

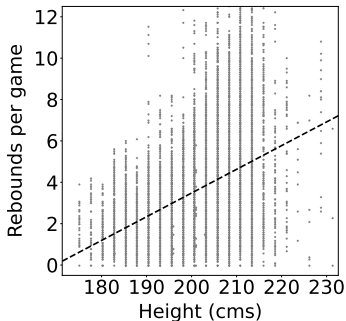
Height and offensive statistics of NBA players between 1996 and 2019

Goal:

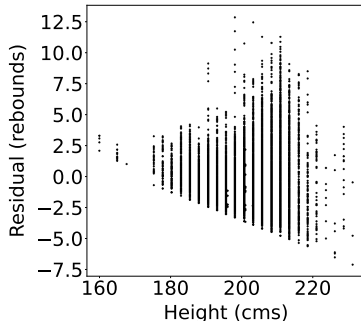
Quantify linear dependence between rebounds/assists/points and height

Rebounds and height

OLS estimator

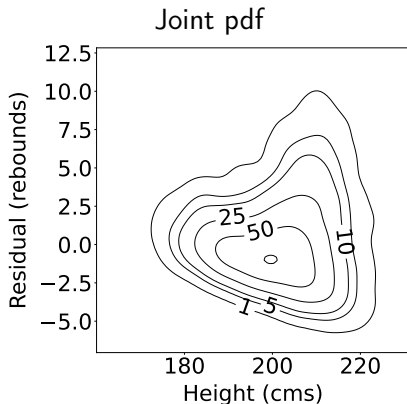
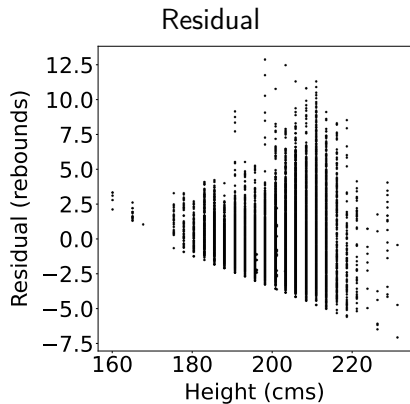


Residual

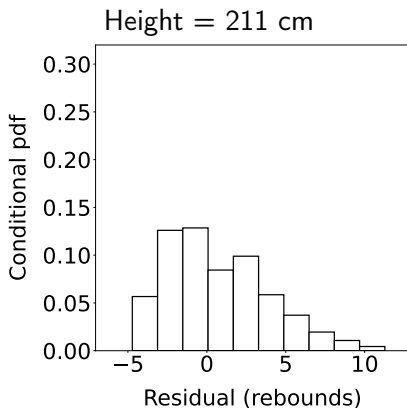
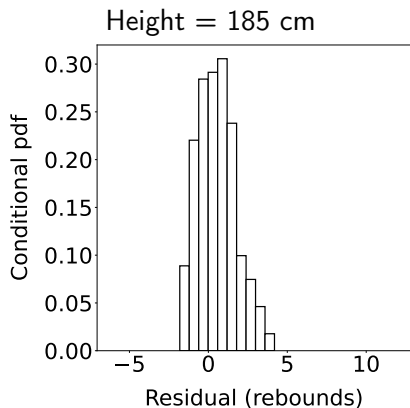


Correlation coefficient between residual and height: 0

Independent?



Conditional distribution of residual given height



What have we learned

- ▶ Independence implies uncorrelation
- ▶ Uncorrelation does not imply independence
- ▶ But it does for Gaussian random variables