

Vorlesung vom 01.11.2012:

- Varianz und Standardabweichung
- Standardisierung als Spezialfall einer linearen Transformation
- Kovarianz und (empirische) Korrelation mehrdimensionaler Merkmale, genauer: von zwei Merkmalen

1) Varianz = mittlere quadratische Abweichung vom arithmetischen Mittelwert

$\vec{x} = \underline{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n = \underbrace{\mathbb{R} \times \dots \times \mathbb{R}}_{n\text{-mal}}$, \mathbb{R} : Menge der reellen Zahlen

Kartesisches Produkt

↑ Datenvektor ↑ Einzeldaten $x_i = X(\omega_i), \omega_i \in \Omega$

Dann, $\text{var}(\vec{x}) = \text{var}(\underline{x}) = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x})^2$ mit

Ist eine Zahl!!

$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ arithm. Mittelwert.

Spezielle Formel für die Varianz (Herleitung):

$$\text{var}(\vec{x}) = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \left\{ \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2 \right\}$$

$= x_i^2 - 2x_i\bar{x} + \bar{x}^2$

$$= (x_1^2 - 2x_1\bar{x} + \bar{x}^2) + (x_2^2 - 2x_2\bar{x} + \bar{x}^2) + \dots + (x_n^2 - 2x_n\bar{x} + \bar{x}^2)$$
$$= (x_1^2 + x_2^2 + \dots + x_n^2) - 2\bar{x}(x_1 + x_2 + \dots + x_n) + (\bar{x}^2 + \bar{x}^2 + \dots + \bar{x}^2)$$

n -mal

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

2) Lineare Transformation und Mittelwert/Varianz:

$y_i = \alpha + \beta x_i$ ($i=1, \dots, n$) Kurz: $\underline{y} = \underline{\alpha} + \beta \underline{x} = \underline{\alpha} + \beta \underline{x}$

neue Daten
Skizze:

alle Daten

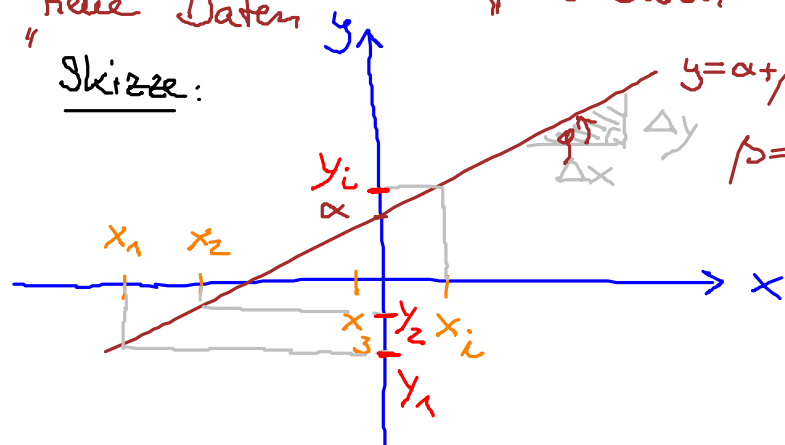
Datenvektor

Datenvektor alt

neue mit $\underline{\alpha} = \underline{\alpha} = \begin{bmatrix} \alpha \\ \alpha \\ \vdots \\ \alpha \end{bmatrix} = \alpha \cdot \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$

$y = \alpha + \beta x$

$\beta = \frac{\Delta y}{\Delta x} = \tan(\varphi)$

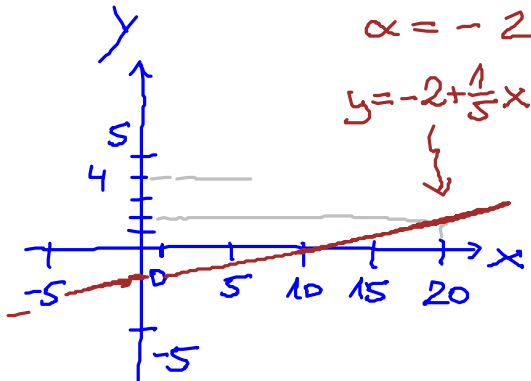


Bsp:

$\underline{x} = (x_1, x_2, x_3, x_4, x_5) = (1, -1, 10, 5, 20)$

$\alpha = -2, \beta = \frac{1}{5} : y = \alpha + \beta x = -2 + \frac{1}{5}x$

$\underline{y} = \underline{\alpha} + \beta \underline{x} = -2 \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \frac{1}{5} \cdot \begin{bmatrix} 1 \\ -1 \\ 10 \\ 5 \\ 20 \end{bmatrix} = \begin{bmatrix} -9/5 \\ -11/5 \\ 0 \\ -1 \\ 2 \end{bmatrix}$



$y_1 = -2 + \frac{1}{5}x_1 = -2 + \frac{1}{5} \cdot 1 = -\frac{9}{5}$

$y_2 = -2 + \frac{1}{5}x_2 = -2 + \frac{1}{5}(-1) = -\frac{11}{5}$

$y_3 = -2 + \frac{1}{5}x_3 = -2 + \frac{1}{5} \cdot 10 = 0$

Satz: Ist $y = \alpha + \beta x$ die zugrunde gelegte lineare Transformation

$y_i = \alpha + \beta x_i$ ($i=1, \dots, n$), so gilt für den Mittelwert und

die Varianz der neuen Daten:

$\bar{y} = \alpha + \beta \bar{x}$ (alter MW) $\text{var}(\underline{y}) = \beta^2 \text{var}(\underline{x})$ (neue Varianz, alte Varianz)

"Beweis":
 1) $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (\alpha + \beta x_i) = \frac{1}{n} \left\{ \alpha \cdot \sum_{i=1}^n 1 + \beta \cdot \sum_{i=1}^n x_i \right\}$
 $= \frac{\alpha \cdot n}{n} + \beta \cdot \frac{1}{n} \cdot \sum_{i=1}^n x_i = \alpha + \beta \bar{x}$
 2) $\text{var}(\underline{y}) = \frac{1}{n} \cdot \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \cdot \sum_{i=1}^n \{ (\alpha + \beta x_i) - (\alpha + \beta \bar{x}) \}^2$
 $= \frac{1}{n} \cdot \sum_{i=1}^n \{ \beta (x_i - \bar{x}) \}^2 = \beta^2 \cdot \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 = \beta^2 \cdot \text{var}(\underline{x})$

Standardisierung von Daten:

$x_i^* := \frac{x_i - \bar{x}}{s} = -\frac{\bar{x}}{s} + \left(\frac{1}{s}\right) \cdot x_i = \alpha + \beta x_i$ mit $\alpha = -\frac{\bar{x}}{s}, \beta = \frac{1}{s}$

mit $s = \sqrt{\text{var}(\underline{x})}$ Standardabweichung!

Dann: $\bar{x}^* = \alpha + \beta \bar{x} = -\frac{\bar{x}}{s} + \frac{1}{s} \bar{x} = 0$
 $\text{var}(\underline{x}^*) = \beta^2 \cdot \text{var}(\underline{x}) = \frac{1}{s^2} \cdot \text{var}(\underline{x}) = \frac{\text{var}(\underline{x})}{\text{var}(\underline{x})} = 1$

3) Kovarianz, Korrelation:

$\underline{X} = (X, Y)$ sei zweidimensionales Merkmal
 z.B. $(X, Y) = (\text{Gewicht}, \text{Körpergröße})$

Ein einzelnes Datum $\underline{X}(\omega) = (X(\omega), Y(\omega))$

Bsp. $\underline{x}_i = (x_i, y_i)$ für ein Individuum x_i

Dann definiert man den Wert

$$\text{Cov}(\underline{x}, \underline{y}) = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})$$

Datenvektor \rightarrow \underline{x}
 Gezeich \rightarrow \underline{y}
 Datenvektor \rightarrow x_i
 Körpergröße \rightarrow y_i

als Kovarianz von X und Y .

Empirisches Korrelationskoeffizient

$$r(\underline{x}, \underline{y}) := \text{Cov}(\underline{x}^*, \underline{y}^*)$$

standardisierte Daten /
Ausprägungen!!

Bemerkung: Verwendung der speziellen linearen Transformation

$$x_i^* = -\frac{\bar{x}}{s_x} + \frac{1}{s_x} x_i \quad \text{und} \quad y_i^* = -\frac{\bar{y}}{s_y} + \frac{1}{s_y} y_i \quad (i=1, \dots, n)$$

Standardabweichung \rightarrow s_x
 bezüglich \rightarrow x
 Merkmal X

Standardabweichung \rightarrow s_y
 bezüglich \rightarrow y
 Merkmal Y

dann erhält man als alternative Formel für den empirischen Korrelationskoeffizienten:

$$\text{Cov}(\underline{x}^*, \underline{y}^*) = \frac{1}{n} \sum_{i=1}^n x_i^* \cdot y_i^* = \frac{1}{n} \sum_{i=1}^n \frac{x_i - \bar{x}}{s_x} \cdot \frac{y_i - \bar{y}}{s_y}$$

$\bar{x}^* = \bar{y}^* = 0$
 $x_i^* = \frac{x_i - \bar{x}}{s_x}$
 $y_i^* = \frac{y_i - \bar{y}}{s_y}$

$$= \frac{1}{s_x \cdot s_y} \cdot \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{\text{Cov}(\underline{x}, \underline{y})}{\sqrt{\text{Var}(x)} \cdot \sqrt{\text{Var}(y)}}$$

$\sqrt{\text{Var}(x)}$
 $\sqrt{\text{Var}(y)}$

ENDE der Vorlesung !!