

Formelsammlung zur multivariaten Statistik (Jens Mehrhoff)

R-Techniken und Q-Techniken

- Datenmatrix: \mathbf{X} , $[n \times m]$
(n Zeilen: Objekte/Personen, m Spalten: Merkmale/Variablen)
- R-Techniken: $\mathbf{R} = \frac{1}{n} \mathbf{Z}' \mathbf{Z}$, $[m \times m]$, $z_{ij} = \frac{x_{ij} - \bar{x}_i}{s_{x_{ij}}}$
(Matrix der Korrelationskoeffizienten zwischen Merkmalen/Variablen)
- Q-Techniken: $\mathbf{Q} = \frac{1}{m} \mathbf{X} \mathbf{X}'$, $[n \times n]$
(Matrix der Korrelationskoeffizienten zwischen Objekten/Personen)

Hotelling's T^2 : nominale Inputvariable, mehrere metrische Outputvariablen, Hypothesentest

- Einstichprobenfall: $T^2 = n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0) \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0)$
- Zweistichprobenfall: $T^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \mathbf{S}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$, $\mathbf{S} = \frac{n_1 - 1}{n - 2} \mathbf{S}_1 + \frac{n_2 - 1}{n - 2} \mathbf{S}_2$

Diskriminanzanalyse ($p = 2, m = 2$): Q-Technik, Dependenzanalyse, Klasseneinteilung

- Elemente: $k = 1, 2, \dots, n_j$ (x_{kji})
- Gruppen: $j = 1, 2, \dots, p$
- Merkmale: $i = 1, 2, \dots, m$
- Diskriminanzfunktion: $y = a_1 x_1 + a_2 x_2$ $\left(x_2 = \frac{y^*}{a_2} - \frac{a_1}{a_2} x_1 \right)$
- Klassifikationsregel: $y \begin{cases} \geq y^* & \text{Objekt in Gruppe 1} \\ < y^* & \text{Objekt in Gruppe 2} \end{cases}$, $y^* = \frac{n_1}{n} \bar{y}_1 + \frac{n_2}{n} \bar{y}_2$
 $\bar{y}_1 = a_1 \bar{x}_{11} + a_2 \bar{x}_{12}$, $\bar{y}_2 = a_1 \bar{x}_{21} + a_2 \bar{x}_{22}$
- Normalgleichungen: $a_1 = \frac{d_1 s_2^2 - d_2 s_{12}}{s_1^2 s_2^2 - s_{12}^2} = \frac{d_1 - r_{12} \frac{d_2}{s_1 s_2}}{1 - r_{12}^2}$, $a_2 = \frac{d_2 s_1^2 - d_1 s_{12}}{s_1^2 s_2^2 - s_{12}^2} = \frac{d_2 - r_{12} \frac{d_1}{s_1 s_2}}{1 - r_{12}^2}$
 $d_1 = \bar{x}_{11} - \bar{x}_{21}$, $d_2 = \bar{x}_{12} - \bar{x}_{22}$
 $s_1^2 = \frac{1}{n} \left[n_1 \sum_{k_1=1}^{n_1} (x_{k_1 11} - \bar{x}_{11})^2 + n_2 \sum_{k_2=1}^{n_2} (x_{k_2 21} - \bar{x}_{21})^2 \right]$
 $s_2^2 = \frac{1}{n} \left[n_1 \sum_{k_1=1}^{n_1} (x_{k_1 12} - \bar{x}_{12})^2 + n_2 \sum_{k_2=1}^{n_2} (x_{k_2 22} - \bar{x}_{22})^2 \right]$
 $s_{12} = \frac{1}{n} \left[n_1 \sum_{k_1=1}^{n_1} (x_{k_1 11} - \bar{x}_{11})(x_{k_1 12} - \bar{x}_{12}) + n_2 \sum_{k_2=1}^{n_2} (x_{k_2 21} - \bar{x}_{21})(x_{k_2 22} - \bar{x}_{22}) \right]$, $r_{12} = \frac{s_{12}}{s_1 s_2}$

Pfadanalyse: Analyse von Kausalstrukturen (nur Falsifikation)

- $$\begin{array}{ccccc}
 & x & \rightarrow & y & \rightarrow & z \\
 \bullet \text{ Gleichungen:} & & & \uparrow & & \uparrow \\
 & & & u & & v \\
 \end{array}$$
- (1) $y = p_{yx}x + u$
 - (2) $z = p_{zy}y + v$
 - Korrelationskoeffizienten
 - $E[(1) \cdot x]: r_{xy} = E(xy) = p_{yx}E(x^2) + E(ux) = p_{yx}$
 - $E[(2) \cdot y]: r_{yz} = E(yz) = p_{zy}E(y^2) + E(vy) = p_{zy}$
 - $E[(2) \cdot x]: r_{xz} = E(xz) = p_{zy}E(xy) + E(vx) = r_{xy}r_{yz}$
 - Bestimmtheit und Unbestimmtheit
 - $E[(1) \cdot u]: E(uy) = p_{yx}E(ux) + E(u^2) = \sigma_u^2 = 1 - r_{xy}^2$
 - $E[(2) \cdot v]: E(vz) = p_{zy}E(vy) + E(v^2) = \sigma_v^2 = 1 - r_{yz}^2$

Faktorenanalyse: R-Technik, Interdependenzanalyse, metrische Inputvariable, latente Outputvariable, Dimensionsreduktion

- Objekte/Personen: $k = 1, 2, \dots, n$
- Merkmale/Variablen: $i = 1, 2, \dots, m$
- Faktoren: $j = 1, 2, \dots, p$
- $\mathbf{Z} = \mathbf{YF}'$
 - $\mathbf{Y} = \mathbf{ZF}(\mathbf{F}'\mathbf{F})^{-1}$, $[n \times p]$: Faktorwerte (unkorrelierte Faktoren)
 - \mathbf{F} , $[m \times p]$: Faktorladungen der Merkmale/Variablen in den Faktoren
- Fundamentaltheorem: $\mathbf{R} = \frac{1}{n}\mathbf{Z}'\mathbf{Z} = \frac{1}{n}\mathbf{FY}'\mathbf{YF}' = \mathbf{FF}' \left(\frac{1}{n}\mathbf{Y}'\mathbf{Y} = \mathbf{I} \right)$
 - \mathbf{R} , $[m \times m]$: Korrelation zwischen Merkmalen/Variablen
- Extraktion der Faktoren: $(\mathbf{R} - \lambda_1\mathbf{I})\mathbf{f}_1 = 0$, $[(\mathbf{R} - \mathbf{f}_1\mathbf{f}_1') - \lambda_2\mathbf{I}]\mathbf{f}_2 = (\mathbf{R}_{(1)} - \lambda_2\mathbf{I})\mathbf{f}_2 = 0, \dots$
- Eigenwertproblem: $|\mathbf{R} - \lambda_1\mathbf{I}| = 0$, $|\mathbf{R}_{(1)} - \lambda_2\mathbf{I}| = 0, \dots$
 - (nicht-triviale Lösungen eines homogenen linearen Gleichungssystems)
 - $\lambda_1, \lambda_2, \dots$: Eigenwerte
 - $\mathbf{f}_1, \mathbf{f}_2, \dots$: Eigenvektoren
- Kommunalitätenproblem: $r_{ii} = 1 - u_{ii} = h_i^2$ (u_{ii} : Varianzanteil des Einzelrestfaktors, h_i^2 : durch die gemeinsamen Faktoren erklärter Anteil der Gesamtvarianz eines Merkmals/einer Variablen)
 - Hauptkomponentenanalyse: $\mathbf{R} = \mathbf{FF}'$ (Zusammenfassung von auf einen Faktor hoch ladenden Merkmalen/Variablen)
 - Faktorenanalyse: $\mathbf{R}^* = \mathbf{R} - \mathbf{U}$ (Ursachenbezeichnung für hohe Ladungen von Merkmalen/Variablen auf einen Faktor)
- Abbruchkriterium: Scree-Test (Anordnung der Eigenwerte der Größe nach)
- Rotationsproblem: $\mathbf{T} = \begin{bmatrix} \cos \delta & -\sin \delta \\ \sin \delta & \cos \delta \end{bmatrix}$ (Drehung gegen den Uhrzeigersinn)
 - $\mathbf{F}^* = \mathbf{FT}$, $\mathbf{R}^* = \mathbf{F}^*\mathbf{F}^{*'} = \mathbf{FTT}'\mathbf{F}' = \mathbf{FF}'$ ($\mathbf{TT}' = \mathbf{I}$)

Latent Structure Analysis (Latent Dichotomy): R-Technik, dichotome Inputvariable, latente Outputvariable, Klasseneinteilung

- relative Klassengrößen: V_I, V_{II}
- bedingte Wahrscheinlichkeiten für positive Ausprägung von i, j : $p_{li}, p_{lli}, p_{lj}, p_{llj}$
(Latent Marginals)
- lokale Unabhängigkeit: $p_{lij} = p_{li}p_{lj}, p_{llij} = p_{lli}p_{llj}, p_{lijk} = p_{li}p_{lj}p_{lk}, p_{llijk} = p_{lli}p_{llj}p_{llk}$
- Accounting Equations
 $V_I + V_{II} = 1, p_i = V_I p_{li} + V_{II} p_{lli}, p_{ij} = V_I p_{lij} + V_{II} p_{llij}, p_{ijk} = V_I p_{lijk} + V_{II} p_{llijk}$
- Kreuzprodukte (Kovarianz): $|ij| = p_{ij} - p_i p_j$
- geschichtete Kreuzprodukte (bedingte Kovarianz): $|ij, k| = p_k p_{ijk} - p_{ik} p_{jk}$
- Steigungen der Trace Lines (Regressionslinien): $b_i = \sqrt{\frac{|ij||ik|}{|jk|}}$
- Hilfsparameter: $\mu_k = p_{lk} p_{llk} = \frac{|ij, k|}{|ij|}$
- quadratische Gleichung: $t^2 + \phi_i t - 1 = 0, \phi_i := \frac{\mu_i}{p_i b_i} - \frac{p_i}{b_i} + \frac{b_i}{p_i}$ (Schiefe)
- Anteilswerte der latenten Klassen: $t = \sqrt{\frac{V_I}{V_{II}}} = \sqrt{\frac{V_I}{1-V_I}}, V_I = \frac{t^2}{1+t^2}, V_{II} = 1 - V_I$
- Position auf der latenten Variablen: $x_I = \frac{1}{t}, x_{II} = -t$
- Latent Marginals: $p_{li} = p_i + b_i x_I = p_i + \frac{b_i}{t}, p_{lli} = p_i + b_i x_{II} = p_i - b_i t$
- Response Pattern (Reproduktion der Daten): $n_I p_{I1} p_{I2} (1 - p_{I3}) = n V_I p_{I1} p_{I2} (1 - p_{I3})$
(Beispiel: ++-, Klasse I)

logistische Regression

- logistische Verteilung: $k = E(y | x) = \beta x_i$
- Wahrscheinlichkeit für positive Ausprägung: $W(X=1) = \frac{\exp(K)}{1 + \exp(K)} = p$
- Wahrscheinlichkeit für negative Ausprägung: $W(X=0) = \frac{1}{1 + \exp(K)} = 1 - p$
- Odd Ratio: $\frac{W(X=1)}{W(X=0)} = \frac{p}{1-p} = \exp(K)$
- Logit: $\ln \frac{p}{1-p} = k$

Data Mining

- Minkowski r -Metriken: $d_r(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_{l=1}^p |x_{il} - x_{jl}|^r \right)^{\frac{1}{r}}$ (nicht skaleninvariant)
- Mahalanobis-Distanz: $d_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j) \mathbf{K}^{-1} (\mathbf{x}_i - \mathbf{x}_j)}, \mathbf{K} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})$