

MVZV $F_X(x_1, \dots, x_p) = P(X_1 \leq x_1, \dots, X_p \leq x_p) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_p} f_X(t_1, \dots, t_p) dt_1 \dots dt_p$ Verteilungsfunktion
 Gemeinsame Dichte von $X \in \mathbb{R}^p$ p-dim ZV

- Randverteilung: $F_{X_i}(x) = P(X_i \leq x) \in \mathbb{R}^p$: (X_1, X_2) $F_{X_1}(x) = P(X_1 \leq x) = \int_{-\infty}^{\infty} \int_{-\infty}^x f(t_1, t_2) dt_1 dt_2$
- Unabhängigkeit: $F_X = \prod_i F_{X_i}(x_i) \forall x_i$
- Bedingt: $f_X(x|y) = \frac{f_X(x,y)}{f_Y(y)}$ stetig, $P(X=x|Y=y) = \frac{P(X=x \cap Y=y)}{P(Y=y)}$ diskret
 $\hookrightarrow f_X(x)$ $\hookrightarrow P(X=x)$
- $E: E(X|Y=y) = \int x f_X(x|y) dx$ stetig, $E(X|Y=y) = \sum_i x_i P(X=x_i|Y=y)$ diskret
 $E(X|Y) \in E(X|Y)(\omega) = E(X|Y=Y(\omega)) \Rightarrow E(E(X|Y)) = E(X)$, $E(BX+b) = BE(X)+b$, $B \in \mathbb{R}^{n \times p}$, $b \in \mathbb{R}^n$, $a \in \mathbb{R}^p$

Varianz $Cov(aX+b, cY+d) = ac Cov(X,Y)$, $Cov(X+Y, Z) = Cov(X,Z) + Cov(Y,Z)$ $|Z|=det(Z)$

$E(X) = (E[X_i])$
 $Cov(X) = \Sigma(X) = (\sigma_{ij})$, $\sigma_{ij} = Cov(X_i, X_j) = E(X_i X_j) - E(X_i)E(X_j)$ sym
 $Cor(X) = \rho(X) = (\rho_{ij})$, $\rho_{ij} = Corr(X_i, X_j) = Cov(X_i, X_j) / (\sigma_{X_i} \sigma_{X_j})$ sym
 $Cov(BX+b) = B \Sigma B^T$, $Var(a^T X) = a^T \Sigma a$
 • Standardisiert: $X^* = \frac{X_i - E(X_i)}{\sigma_{X_i}} = D^{-1}(X - E(X))$, $D = eye(\sigma_i)$
 • Unabhängig $X \sim N(\mu, \Sigma)$: X_i u $\Leftrightarrow \Sigma$ diag $\Leftrightarrow X_i$ unkorreliert
 • Schätzer: $\bar{X} = \frac{1}{n} \sum X_i$, $S = \frac{1}{n-1} \sum (X_i - \bar{X})(X_i - \bar{X})^T = (s_{ij})$, $R = (r_{ij}) = (s_{ij} / \sqrt{s_{ii} s_{jj}})$

Normalverteilung
 $f_X(x) = (\sigma \sqrt{2\pi})^{-p/2} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$
 $f_X(x) = (2\pi)^{-p/2} |Z|^{-0.5} \exp(-\frac{1}{2}(x-\mu)^T Z^{-1} (x-\mu))$
 $\hookrightarrow X_i$ u stdnorm V $\Rightarrow AX+b$ multi norm V mit $\mu = b$, $Z = A^T A$

Verlustfunktion: $L(y, \hat{f}(x))$, Trainingsfehler: $\bar{err}(L) = \frac{1}{n} \sum L(y_i, \hat{f}(x_i))$ Bsp: $L = (y - \hat{f}(x))^2$
 Testfehler: $E(L) = \int L(y, \hat{f}(x)) f_{X,Y}(x,y) dx dy$

Diskret Unabhängig

	b_1	a_1	c_1	gew
Y	b_n	a_n	c_n	1
	p_1	p_1	p_1	p_1

Multi Lin Reg
 • Modell: $Y = f(X) + \epsilon = Xb + \epsilon \Leftrightarrow \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} b_0 \\ \vdots \\ b_p \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$
 $\begin{cases} Cov(\epsilon_i, \epsilon_j) = 0 \\ E(\epsilon_j) = 0 \\ Var(\epsilon_j) = \sigma^2 \end{cases}$ Erklärte Varianz $V(Y)/V(X) \leftarrow Var(Variable)/Var(Gesamt$
 $\hookrightarrow \hat{b} = (X^T X)^{-1} X^T Y \hookrightarrow E(Y|X=x) = b_0 + \sum_{j=1}^p x_{j1} b_{j1} \leftarrow \hat{b}_0 + \sum_{j=1}^p x_{j1} \hat{b}_{j1} = (1, x^T) \hat{b}$ ist Schätzer
 $\hookrightarrow \hat{y} = X \hat{b}$ geschätzte Werte $\hookrightarrow \hat{\beta} = Hy$ mit $H = X(X^T X)^{-1} X^T$ sym hat Matrix
 $\hookrightarrow \hat{e} = y - \hat{y}$ Residuenvektor $\Rightarrow SSE = \hat{e}^T \hat{e} = \sum (y_i - \hat{y}_i)^2$ Residual sum of squares \leftarrow wird durch CS min
 $\hookrightarrow \hat{\sigma}^2 = \frac{1}{n-(p+1)} \hat{e}^T \hat{e} = \frac{1}{n-(p+1)} \sum (\hat{e}_i - \bar{\hat{e}})^2$ Residual variance $\Rightarrow RSE = \sqrt{\hat{\sigma}^2}$ Residual standard error
 $\hookrightarrow E(\hat{b}) = b$, $Cov(\hat{b}) = \sigma^2 (X^T X)^{-1}$, $E(\hat{\sigma}^2) = \sigma^2$, $E(\hat{e}) = 0 \leftarrow \hat{\sigma}^2$ Approximation: $Cov(\hat{b}) = \hat{\sigma}^2 (X^T X)^{-1}$
 • VIF = $1/(1-R_j^2)$ Steigerungsfaktor von b_j im Vergleich zu X_j unkorreliert \leftarrow < 4 unbedenklich
 • Statistische Tests. Annahme: $\epsilon \sim N(0, \sigma^2 I_n)$ verteilt
 $\hookrightarrow \hat{b} \sim N_{p+1}(b, \sigma^2 (X^T X)^{-1}) \rightarrow \hat{b}_j - b_j / \hat{\sigma} \sqrt{v_j} \sim N(0, 1)$ mit $v_j = \text{spur}[(X^T X)^{-1}]_j$
 $\hookrightarrow CI: [\hat{b}_j \pm z_{1-\alpha/2} \hat{\sigma} \sqrt{v_j}] \hookrightarrow H_0: \hat{b}_j = 0 \rightarrow \hat{b}_j / \hat{\sigma} \sqrt{v_j} \sim t_{n-(p+1)}$ R^2 Anteil der Variabilität der Y_i die durch das Modell erklärt wird

Anpassungsgrade
 $\hookrightarrow SST = \sum (y_i - \bar{y})^2$ total sum of squares
 $\hookrightarrow SSE = \sum (y_i - \hat{y}_i)^2$ residual sum of squares
 $\hookrightarrow SS M = \sum (\hat{y}_i - \bar{y})^2$ model sum of squares
 $R^2 = \frac{SSM}{SST} = \frac{SST - SSE}{SST}$ $SST = SSE + SSM$
 $R^2 = Cor(Y, \hat{Y})^2$
 $p=1 \rightarrow Cor(Y, X)^2$

Model parameter [p Kovariaten $\Rightarrow p+1$ Modellparameter]
 $\hookrightarrow AIC = n \log(\frac{1}{n} SSE) + 2(p+1)$
 $\hookrightarrow BIC = n \log(\frac{1}{n} SSE) + \log(n)(p+1)$ aus Modell mit allen Kovariaten
 $\hookrightarrow \hat{\sigma}^2 = SSE / (n - (p+1))$, $\hat{\sigma}^2$ Varianzschätzer

Variablen Selektion \leftarrow Subset: alle, berechnen
 Forward: $p=0$ und dann nach best hinzufügen \leftarrow backward: $p=\text{voll}$ nach best entfernen

$\hookrightarrow Adj R^2 = 1 - \frac{SSE/(n-p-1)}{SST/(n-1)}$ \leftarrow X_1, \dots, X_p $q < p$
 $\hookrightarrow F = \frac{(SSE^* - SSE)/(p-q)}{SSE/(n-p-1)} \sim F_{p-q, n-p-1}$ unter $H_0: \beta_{q+1} = \dots = \beta_p = 0$

Klassifikation $X: (\mathbb{R}^1 \times \mathbb{R}^p) \rightarrow \{1, \dots, K\}$, $g(X) = \begin{cases} 0, & b^T x \leq 0 \\ 1, & b^T x > 0 \end{cases}$ Accuracy = $(TP+TN)/Total$
 • Konfusions Matrix $y = h(x) = \sum_{i=0}^p b_i x_i \in b_i$ Diskriminanzkoeffizient
 • Missclassification Rate = $(FP+FN)/Total$

	$Y=1$	$Y=0$
$Y=1$	TP	FP
$Y=0$	FN	TN

 TP = True positive, FP = False positiv, Fehler 1. Art
 FN = False negativ, Fehler 2. Art

Lin Diskriminanz $[X|Y=k \sim N(\mu_k, \Sigma)] \forall k=1, \dots, K$
 \hookrightarrow Annahme: Gemeinsame Kovarianz, fkt Dichten der multivariaten Normalverteilung
 $\hookrightarrow \pi_k = n_k/n$
 $\hookrightarrow K=2 \leftarrow g(x) = 1 \Leftrightarrow h(x) > 0$
 $h(x) = x^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_0) - \frac{1}{2} (\hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 - \hat{\mu}_0^T \hat{\Sigma}^{-1} \hat{\mu}_0) + \log(\frac{\hat{\pi}_1}{\hat{\pi}_0})$
 \hookrightarrow Quadratisch: $h(x) = \log(\pi_1) - \log(\pi_0) = -\frac{1}{2} \log(|\Sigma_1|) - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) - \log(\hat{\pi}_1) + \log(\hat{\pi}_0)$

(Bsp) $\pi_0 = 0.5, \pi_1 = 0.5$ $\pi_0 = 0.3, \pi_1 = 0.7$

Multitest $P = P_{H_0}(T(X_1, \dots, X_n) > T(x_1, \dots, x_n))$
 $P \leq \alpha \Leftrightarrow T(x_1, \dots, x_n) > c_\alpha \Leftrightarrow H_0 \text{ verwerfen}$

# wahre H_0	# Ablehnung	# Nicht	
V	V	$m_0 - V$	m_0
# falsche H_0	$R - V$	$m_1 - R + V$	m_1
	R	$m - R$	m

Linke
 $C(C_i, C_j)$
 (cluster Abstand)
 (single linkage)
 (complete linkage)
 (centroid)

$\square \text{FWER} = P(V > 0)$ $\square \text{FDR} = E(V / \max(R, 1))$ \leftarrow p-werte unabhängig
 Bonferroni: $P_i \leq \frac{\alpha}{m}$; Sidak: $P_i \leq 1 - (1 - \alpha)^{\frac{1}{m}}$; Holm: $P_i \leq \frac{\alpha}{m}$, p-Werte sortiert, $m = m - 1$

cluster $d(i, j) = d(j, i)$; $d(i, j) \geq 0$; $d(i, i) = 0$; Agglomerativ: n cluster \rightarrow nearest merge
 k-means: k cluster \rightarrow d min zuweisen \rightarrow C_k Mittelpunkt neue $C_k \rightarrow$ repeat bis no change
 Single: $\min_{x \in Y, y \in Y} d(x, y)$
 comp: $\max_{x \in Y, y \in Y} d(x, y)$

PCA $\Sigma = \text{Cov}(X) = A \Lambda A^T$, A normierte EV Ladungsmatrix, $\Lambda = \text{diag}(\lambda_i)$ $\lambda_1 \geq \dots \geq \lambda_p \geq 0$
 $Y = A^T X \Rightarrow Y = A^T X$, $Y = A_i^T X$ i-te Hauptkomponente zu X
 $\text{Cov}(Y_i, Y_j) = 0$, $\text{Var}(Y_i) = \lambda_i$, $\sum \text{Var}(Y_i) = \text{spur}(\Sigma) = S$ Gesamtvarianz $\Rightarrow \lambda_i$ erklärte Varianz
 Kaiser-Kriterium: $\lambda_i \geq \frac{1}{p} \sum \lambda_i$

Multi Lin Reg Überprüfung der Modellvoraussetzungen

Modellannahme	Konsequenzen aus Verletzung	Überprüfung
A Ein Zusammenhang zwischen X, Y	Verzerrung der Schätzer	Residuenplot
B $E(\epsilon_i) = 0$, $\text{Var}(\epsilon_i) = \sigma^2$	Verzerrung der Schätzer	Residuenplots (\hat{y} vs $\hat{\epsilon}^*$, x_j vs $\hat{\epsilon}^*$)
C Unkorreliertheit der Fehlerterme	Verzerrung der Schätzer	zeitliche Strukturen
D Normalverteilung der Fehlerterme	stat. Tests und CI nicht auf Niveau	QQ-Plots der Residuen
E $\text{Rang}(X) = p+1 \Rightarrow X_j$ lin unabh	Parameterschätzer nicht eindeutig	Korrelationsplots
F Multikollinearität	instabile Schätzer, hohe SE	Korrelationsplots, VIF

$\square A+B+C \Rightarrow$ best linear unbiased estimator \Rightarrow unverzerrte Schätzer
 \square Residuenplots sollten keine Struktur aufweisen
 \square Verletzungen in A, B, D können durch Transformation behoben werden
 \square Verletzung in F können Kovariaten entfernt werden (PCA)

Shapiro-Wilk Test auf Normalverteilung
 $p > \alpha \rightarrow$ Normalverteilung
 Log-Likelihood \leftarrow Variable aufnehmen, wenn $T > \chi^2_{1-\alpha}$
 $T = 2 \ln \frac{L(\mu_0)}{L(\mu_1)} \sim \chi^2_{1-\alpha}$

R-Output • Multiple R^2 : Anpassungsgüte. Wie nahe die Punkte der Geraden sind. Beschreibt die erklärte Varianz und steigt mit der Anzahl an Kovariaten.
 • Adjusted R^2 : R^2 von allen statistisch signifikanten Kovariaten
 • F-Statistic: model MSE/residual MSE. Test der Nullhypothese, dass alle Parameter 0 sind.

ANOVA	DF	sum squares	mean sum squares	F-Test
Regression (Modell)	p	SSM	SSM/p	$F = \frac{SSM/p}{SSE/[n-(p+1)]}$
Residuen (Error)	n-(p+1)	SSE	SSE/[n-(p+1)]	
Gesamt (Total)	n-1	SST	SST/(n-1)	

Variablen Selektion [Best: min AIC, max $R^2 \Leftrightarrow$ min SSE]
 Subset selection: Alle Modelle berechnen \Rightarrow choose best
 Forward Selection: Starte mit p=0 \Rightarrow choose best new Kovariate

Diskriminanz

\square Bayes-classifier: $\pi_k(x) = P(X=k|X=x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l} \Rightarrow \hat{y} = \arg \max_k \pi_k(x)$
 $\leftarrow \hat{\pi}_k(x)$ Schätzer durch $f_k(x) = \hat{f}_k(x) \cdot (I|X=l)$
 \square Logistische Regression [Y binär {0,1}, X p-dim, $x_i \in \mathbb{R}^p$]
 $\hookrightarrow P(Y=1|X=x) = \frac{\exp(b^T x)}{1 + \exp(b^T x)}$, $P(Y=0|X=x) = \frac{1}{1 + \exp(b^T x)}$
 $\hookrightarrow L(b) = \prod_{i=1}^n P_b(Y_i = k_i | X = x_i) \hookrightarrow \text{odds}(x) = \exp(b^T x) \leftarrow \exp(b_i)$ relative Veränderung der odds bei Veränderung von Merkmal i
 \square Diskriminierung
 \hookrightarrow Güte anhand der Trennschärfe
 \hookrightarrow Fehlerrate: $P(Y \neq \hat{Y}) \hookrightarrow$ Sensitivität: $P(\hat{Y}=1|Y=1)$ \hookrightarrow Spezifität: $P(\hat{Y}=0|Y=0)$
 \hookrightarrow positiver prädiktiver Wert: $P(Y=1|\hat{Y}=1)$ \hookrightarrow negativer prädiktiver Wert: $P(Y=0|\hat{Y}=0)$

