

# Histogramm

$$\text{Histogramm-Funktion: } H(x) = \sum_{i=1}^{k-1} \left( \frac{n_i}{n} \frac{1}{t_{i+1} - t_i} \right) I_{[t_i, t_{i+1})}(x)$$

k... Klassenanzahl, t ... Intervallgrenzen I ... Indikatorfunktion, n... Umfang Stichprobe,  $n_i$ ... # im Intervall

$$\text{Intervalllänge nach Sturges: } h_n = t_{i+1} - t_i = \lceil \log_2(n) + 1 \rceil$$

$$\text{Intervalllänge nach Scott: } h_n = \frac{3.5s}{\sqrt[3]{n}} \quad s \dots \text{ empirische Standardabweichung}$$

$$\text{Intervalllänge nach Freedman: } h_n = \frac{(2IQR)}{\sqrt[3]{n}}$$

## Verteilungsfunktion

von 0 bis 1, monoton wachsend,  $\lim_{x \rightarrow \infty} F(x) = 1$ ,  $\lim_{x \rightarrow -\infty} F(x) = 0$

$$\text{theoretische Verteilungsfunktion: } F(x) = \int_{-\infty}^x f(t) dt$$

$$\text{empirische Verteilungsfunktion: } F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[x_i, \infty)}(x)$$

# Dichteschätzung/Kernschätzung

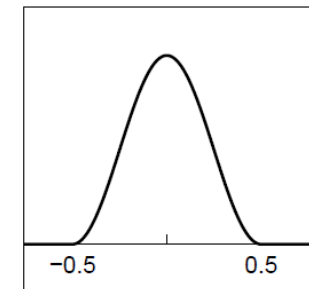
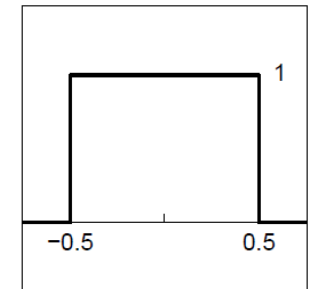
$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} W\left(\frac{x - x_i}{h}\right)$$

$W(t)$  ... Gewichtsfunktion  $\int_{-\infty}^{\infty} W(t) = 1$ ,

$\left[x - \frac{h}{2}, x + \frac{h}{2}\right]$  ... Fenster ( $h$  = Intervalllänge)

Rechtecks Gewichtsfunktion:  $W(t) = \begin{cases} 1 & |t| \leq \frac{1}{2} \\ 0 & \text{sonst} \end{cases}$

Cosinus Gewichtsfunktion:  $W(t) = \begin{cases} 1 + \cos 2\pi t & |t| < \frac{1}{2} \\ 0 & \text{sonst} \end{cases}$



## Wahrscheinlichkeitsnetz

=Verteilungsfunktion mit verzerrter vertikaler Achse  $\rightarrow$  zb.: Abstände proportional zur Inversen der Verteilungsfunktion von  $N(0,1) = \Phi$ . y-Achse =  $\Phi^{-1}(F_n(x))$  x-Achse =  $x$ .

Bei ungefähr Normalverteilung gilt:  $\Phi^{-1}(F_n(x)) \sim \Phi^{-1}\left(\Phi\left(\frac{x-\mu}{\sigma}\right)\right) = \frac{x-\mu}{\sigma} \rightarrow$  auf einer Geraden.

bei Wahrscheinlichkeit 50% Schätzung für  $\mu$  bei 84% Schätzung für  $\mu + \sigma$

# Quantile-Quantile Plots

$$F_x(t) = P(X \leq t), q_x(p) = F_x^{-1}(p)$$

$F_x = F_y \Leftrightarrow (q_x(p_i), q_y(p_i))$  liegen auf der Geraden  $y = x$

## univariate Schätzer

arithmetisches Mittel:  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  empirische Standardabweichung  $s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

gestutztes Mittel:  $m(\alpha) = \frac{1}{n-2g} (x_{(g+1)} + \dots + x_{(n-g)})$   $\alpha \dots 0 \leq \alpha \leq 0.5$   $g = \lfloor n\alpha \rfloor$  (auf

Ganzzahl abgerundet)

gestutzte Streuung  $S(\alpha) = \sqrt{\frac{1}{n-2g-1} \sum_{i=g+1}^{n-g} (x_{(i)} - m(\alpha))^2}$

gestutzte SD =  $\frac{S(\alpha)}{c_\alpha}$   $c_\alpha \dots$  abhängig von  $\alpha$

$$\text{Median } \mathbf{median}(x_1, \dots, x_n) = \begin{cases} \mathbf{x}_{\left(\frac{n+1}{2}\right)} & n \text{ ungerade} \\ \frac{\mathbf{x}_{\left(\frac{n}{2}\right)} + \mathbf{x}_{\left(\frac{n}{2}+1\right)}}{2} & n \text{ gerade} \end{cases}$$

$$\text{MAD (Median Absolute Deviation)} = \mathbf{median}_{1 \leq i \leq n} \left( \left| \mathbf{x}_i - \mathbf{median}_{1 \leq j \leq n}(\mathbf{x}_j) \right| \right)$$

$$s_{MAD} = \frac{MAD}{0.675} = 1.483 \cdot MAD$$

$$\text{25 Quartil } \mathbf{Q}_{0.25} = \mathbf{median} \left( \mathbf{x}_{(1)}, \dots, \mathbf{x}_{\left(\left[\frac{n+1}{2}\right]\right)} \right)$$

$$\text{75 Quartil } \mathbf{Q}_{0.75} = \mathbf{median} \left( \mathbf{x}_{\left(\left[\frac{n}{2}+1\right]\right)}, \dots, \mathbf{x}_{(n)} \right)$$

$$\text{Interquartilabstand } \mathbf{IQR} = \mathbf{Q}_{0.75} - \mathbf{Q}_{0.25} \quad s_{IQR} = \frac{IQR}{1.35}$$

$$Q_n \text{ Streuungsschätzer } \mathbf{Q}_n = \left\{ \left| \mathbf{x}_i - \mathbf{x}_j \right|; i < j \right\}_k \quad s_{Q_n} = 2.219 \cdot Q_n$$

## Dichteschätzung in zwei Dimensionen

$$\text{Boxcar Funktion } W(u, v) = \begin{cases} \frac{1}{\pi} & u^2 + v^2 \leq 1 \\ 0 & \text{sonst} \end{cases}$$

$$\text{Cosinus Gewichtsfunktion } W(u, v) = \begin{cases} \frac{1 + \cos(\pi\sqrt{u^2 + v^2})}{\pi} & \text{wenn } u^2 + v^2 \leq 1 \\ 0 & \text{sonst} \end{cases}$$

$$\text{Dichtefunktion } \hat{f}(x, y) = \frac{1}{h^2 n} \sum_{i=1}^n W\left(\frac{x-x_i}{h}, \frac{y-y_i}{h}\right) \quad h \dots \text{ Fensterbreite}$$

## Robuste Schätzung linearer Trends

$$\mathbf{y} = \alpha + \beta \mathbf{x} + \boldsymbol{\varepsilon} \quad \alpha \dots \text{Abszissenabstand} \quad \beta \dots \text{Steigung} \quad \boldsymbol{\varepsilon} \dots \text{Fehler}$$

$$\mathbf{r}_i = \mathbf{y}_i - \hat{\mathbf{y}}_i = \mathbf{y}_i - \hat{\alpha} - \hat{\beta} \mathbf{x}_i \quad i = 1, \dots, n$$

Least Squares

$$(\hat{\alpha}_{LS}, \hat{\beta}_{LS}) = \operatorname{argmin}_{\alpha, \beta} \sum_{i=1}^n (\mathbf{y}_i - \alpha - \beta \mathbf{x}_i)^2 = \operatorname{argmin}_{\alpha, \beta} \sum_{i=1}^n r_i^2$$

$$\hat{\beta}_{LS} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \hat{\alpha}_{LS} = \bar{y} - \hat{\beta}_{LS} \bar{x}$$

Bruchpunkt:  $\frac{k}{n}$  n ... Stichprobenumfang, k ... maximale Anzahl von Ausreißern

Tukey

Datenpaare in 3 Gruppen nach x-Werte:  $n_L + n_M + n_R = n$

Median pro Gruppe  $x_L = \mathit{median}_{(x_i, y_i) \in L} x_i$   $y_L = \mathit{median}_{(x_i, y_i) \in L} y_i$

Schätzwerte  $\hat{\beta}_0 = \frac{y_R - y_L}{x_R - x_L}$ , Residuen  $r_i^0 = y_i - (\hat{\alpha}_0^{(*)} + \hat{\beta}_0(x_i - x_M))$  Bruchpunkt =  $\frac{1}{6}$  (Median eines  $\frac{1}{3}$ )

Theil

$\hat{\beta}_{ij} = \frac{y_j - y_i}{x_j - x_i}$   $1 \leq i < j \leq n$   $\hat{\beta}_T = \mathit{median}_{1 \leq i < j \leq n} (\hat{\beta}_{ij})$  Bruchpunkt = 0.29

Siegel (Repeated Median Line) (Bruchpunkt 0.5)

$\hat{\beta}_{RM} = \mathit{median}_{1 \leq i \leq n} (\mathit{median}_{1 \leq j \leq n} (\hat{\beta}_{ij}))$   $\hat{\alpha}_{RM} = \mathit{median}_{1 \leq i \leq n} (y_i - \hat{\beta}_{RM} x_i)$

LMS (Least Median of Squares) Regression (Bruchpunkt 0.5)

$(\hat{\alpha}_{LMS}, \hat{\beta}_{LMS}) = \mathit{argmin}_{\alpha, \beta} \mathit{median}_i r_i^2$

LTS (Least Trimmed Squares) Regression

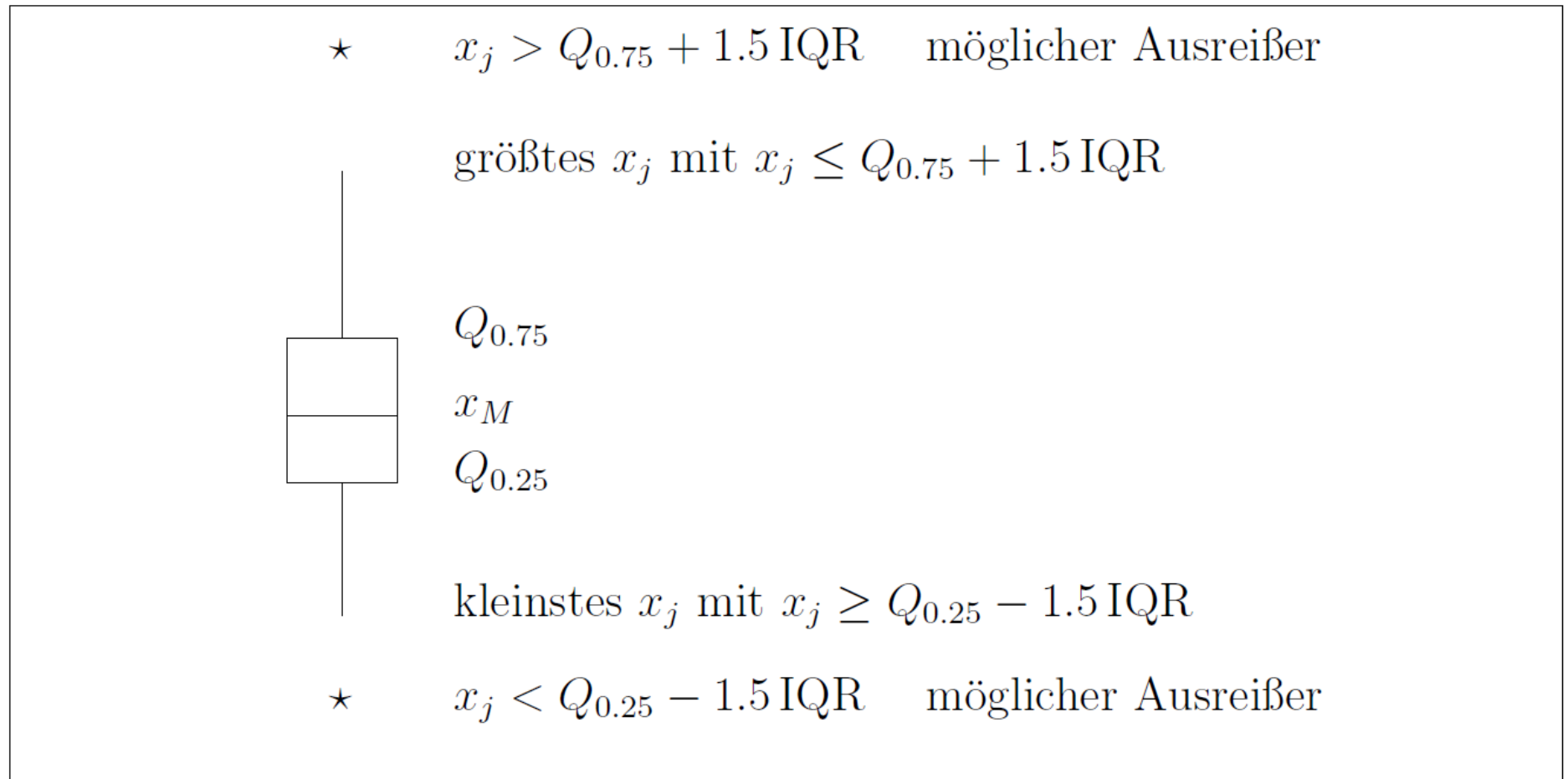
$(\hat{\alpha}_{LTS}, \hat{\beta}_{LTS}) = \mathit{argmin}_{\alpha, \beta} \sum_{i=1}^h r_{(i)}^2$   $\frac{n}{2} < h < n$

Stichprobe:  $x_1, \dots, x_n$

$x_M := \text{median}(x_1, \dots, x_n)$

$Q_{0.25}$  und  $Q_{0.75}$  sind die Quantile 0.25 bzw. 0.75

$\text{IQR} = Q_{0.75} - Q_{0.25}$  ist der Interquartil-Abstand



# Nichtlineare Glättung

Lineare Filter  $\mathbf{z}_t = \sum_{i=-l_1}^{l_2} \alpha_i \mathbf{x}_{t+i}$  gewichtete Summe der  $x_i$  in einer Umgebung von  $t$

Medianglättung  $G\mathbf{x}_t = \mathit{median}(\mathbf{x}_{t-s}, \mathbf{x}_{t-s+1}, \dots)$  Wert an Stelle  $t$  ist Median aus einer  $2s + 1$  Umgebung

Repeated Median

lokale lineare Approximation  $\mu_{t+1} \approx \mu_t + \beta_t \cdot i$   $\mu_t$  ... Level  $\beta_t$  ... Steigung,  $i = -s$  bis  $s$

$$\mathbf{r}_t(t + i) = \mathbf{x}_{t+i} - \hat{\mu}_t - \hat{\beta}_t \cdot i$$

$$\hat{\beta}_t = \mathit{median}_{-s \leq i \leq s} \left( \mathit{median}_{-s \leq j \leq s} \left( \frac{x_{t+i} - x_{t+j}}{i - j} \right) \right) \quad \hat{\mu}_t = \mathit{median}_{-s \leq i \leq s} (\mathbf{x}_{t+i} - \hat{\beta}_t \cdot i)$$

Ausreißer  $|\hat{\mathbf{r}}_t| > 2 \cdot \hat{\sigma}_t$   $\hat{\mathbf{r}}_t = \mathbf{x}_t - \hat{\mu}_t$

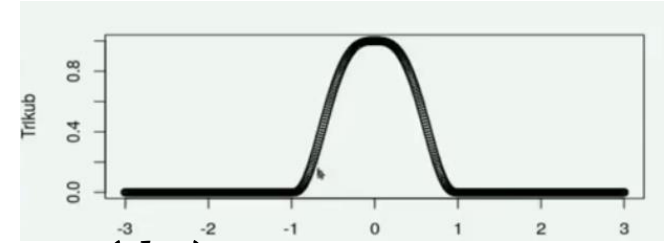


# LOWESS (LOcally WEighted regression Scatter plot Smoothing)

$q = \lfloor nf + 0.5 \rfloor$  ... Anzahl der Punkte die zur Glättung verwendet werden

$d_{ik} = |x_i - x_k|$        $d_i = |x_i - x_{i_{max}}|$        $i_{max}$  ... Index des am weitesten entfernten Punkt aus  $q$

Trikubische Gewichtungsfunktion  $T(t) = \begin{cases} (1 - |t|^3)^3 & |t| < 1 \\ 0 & \text{sonst} \end{cases}$



Gewicht von  $(x_k, y_k)$  bezüglich  $(x_i, y_i)$        $t_i(x_k) = \begin{cases} T\left(\frac{|x_i - x_k|}{d_i}\right) = T\left(\frac{d_{ik}}{d_i}\right) & d_i \neq 0 \\ 1 & d_i = 0 \end{cases}$

gewichtete Regression (Zeitpunkt  $i$ )  $\min \sum_{k=1}^n t_i(x_k) (y_k - a - bx_k)^2$  Residuen  $r_i = y_i - \hat{y}_i$

Biweight Gewichtungsfunktion  $B(t) = \begin{cases} (1 - t^2)^2 & |t| < 1 \\ 0 & \text{sonst} \end{cases}$

$m = \text{median}_{1 \leq k \leq n} |r_k| \rightarrow 3m \approx 2\sigma$  Gewichte Residuen  $w(x_k) = B\left(\frac{r_k}{3m}\right)$

robust gewichtete Regression  $\min \sum_{k=1}^n w(x_k) t_i(x_k) (y_k - a - bx_k)^2$

Upper and Lower Smoothing Aufteilung der Residuen in positive & negative Residuen  $r_i = y_i - \hat{y}_i$

$r_i^+$  für  $(x_i^+, \hat{y}_i^+)$

# Zeitreihenanalyse

$$\mathbf{x}_t = \boldsymbol{\tau}_t + \boldsymbol{\delta}_t + \mathbf{e}_t \quad \tau_t \dots \text{Trendkomponente}, \delta_t \dots \text{Saisonkomponente}, e_t \dots \text{Restkomponente}$$

Lineares Modell  $\ln(\mathbf{x}_t) = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 t + \mathbf{e}_t$

quadratisches Modell  $\ln(\mathbf{x}_t) = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 t + \boldsymbol{\beta}_2 t^2 + \mathbf{e}_t \rightarrow \hat{\mathbf{x}}_t = \exp(\hat{\boldsymbol{\beta}}_0 + \hat{\boldsymbol{\beta}}_1 t + \hat{\boldsymbol{\beta}}_2 t^2)$

Fourierreihe 1. Ordnung  $f(t) = a_0 + \sum_{j=1}^J (a_j \cos(\omega_j t) + b_j \sin(\omega_j t)) \quad \omega_j = \frac{2\pi j}{P}$

$$\ln(\mathbf{x}_t) = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 t + \boldsymbol{\beta}_2 t^2 + \boldsymbol{\beta}_3 \cos\left(\frac{2\pi}{P} \cdot t\right) + \boldsymbol{\beta}_4 \sin\left(\frac{2\pi}{P} \cdot t\right) + \mathbf{e}_t$$

## Exponentielles Glätten

$$\tilde{\mathbf{x}}_t = \alpha \mathbf{x}_t + (1 - \alpha) \tilde{\mathbf{x}}_{t-1} \quad \alpha \dots \text{Glättungsfaktor } 0 < \alpha < 1 \quad \text{Startwert } \tilde{\mathbf{x}}_m = \text{mean}(x_0 \text{ bis } x_m)$$

Prognose  $\tilde{\mathbf{x}}_{t+h|t} = \tilde{\mathbf{x}}_t$

## Glättung nach Holt-Winters

$$\tilde{\mathbf{x}}_t = \alpha \mathbf{x}_t + (1 - \alpha)(\tilde{\mathbf{x}}_{t-1} + \mathbf{b}_{t-1}) \quad \mathbf{b}_t = \beta(\tilde{\mathbf{x}}_t - \tilde{\mathbf{x}}_{t-1}) + (1 - \beta)\mathbf{b}_{t-1}$$

Prognose  $\tilde{\mathbf{x}}_{t+h|t} = \tilde{\mathbf{x}}_t + h\mathbf{b}_t$

Autokovarianz  $\mathbf{Cov}(\mathbf{x}_t, \mathbf{x}_{t-k})$  Schätzer  $\mathbf{c}_k = \frac{1}{T} \sum_{t=k+1}^T (\mathbf{x}_t - \bar{\mathbf{x}})(\mathbf{x}_{t-k} - \bar{\mathbf{x}})$   $k \dots$  lag (Abstand)

Autokorrelation der Ordnung  $k$   $\rho_k = \mathbf{Corr}(\mathbf{x}_t, \mathbf{x}_{t-k}) = \frac{\mathbf{Cov}(\mathbf{x}_t, \mathbf{x}_{t-k})}{\mathbf{Var}(\mathbf{x}_t)}$   $\mathbf{r}_k = \frac{\mathbf{c}_k}{\mathbf{c}_0}$   $\mathbf{c}_0 \dots$  Varianz von  $\mathbf{x}_t$

## Zeitreihenmodelle

Moving Average (MA) Modell MA(1):  $\mathbf{x}_t = \mathbf{a} + \mathbf{u}_t - \Theta_1 \mathbf{u}_{t-1}$   $\mathbf{u}_t \dots$  white noise

Autoregressives (AR) Modell AR(1):  $\mathbf{x}_t = \mathbf{a} + \phi \mathbf{x}_{t-1} + \mathbf{u}_t$   $\mathbf{u}_t \dots$  white noise

ARMA (Autoregressive-Moving-Average) Modell:

ARMA(1,1) =  $\mathbf{x}_t = \mathbf{a} + \phi \mathbf{x}_{t-1} + \mathbf{u}_t - \Theta_1 \mathbf{u}_{t-1}$

ARIMA Modell

Differenz Operator  $\Delta \mathbf{x}_t = \mathbf{x}_t - \mathbf{x}_{t-1}$

$$\Delta^2 \mathbf{x}_t = \Delta(\Delta \mathbf{x}_t) = \Delta(\mathbf{x}_t - \mathbf{x}_{t-1}) = \mathbf{x}_t - 2\mathbf{x}_{t-1} + \mathbf{x}_{t-2}$$

ARIMA(1,1,1):  $\mathbf{x}_t = \mathbf{x}_{t-1} + \mathbf{a} + \phi(\mathbf{x}_{t-1} - \mathbf{x}_{t-2}) + \mathbf{u}_t - \Theta_1 \mathbf{u}_{t-1}$ ,

## Kovarianz

$$\sigma_{jk} = E \left[ \left( x_j - E(x_j) \right) \left( x_k - E(x_k) \right) \right], \quad (j, k \in \{1, \dots, p\} \text{ } p \dots \# \text{Variablen})$$

$$\text{für } j = k: \sigma_{jj} = E \left[ \left( x_j - E(x_j) \right)^2 \right] = \text{Varianz}$$

## Stichprobenkovarianz

$$s_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

## Korrelationskoeffizient

$$\rho_{jk} = \frac{\sigma_{jk}}{\sqrt{\sigma_{jj}\sigma_{kk}}} = \frac{\text{Kovarianz } jk}{\sqrt{\text{Varianz}(j) \cdot \text{Varianz}(k)}}$$

## Stichprobenkorrelation

$$r_{jk} = \frac{s_{jk}}{\sqrt{s_{jj}s_{kk}}}$$

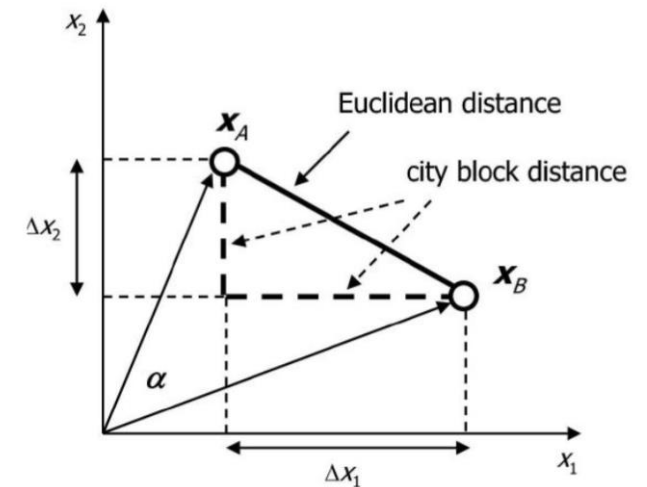
# Distanzmaß

Euklidische Distanz  $\mathbf{d}_E(\mathbf{x}_A, \mathbf{x}_B) = \left( \sum_{j=1}^p (\mathbf{x}_{Bj} - \mathbf{x}_{Aj})^2 \right)^{\frac{1}{2}} = \|\mathbf{x}_B - \mathbf{x}_A\|$

Manhattan Distanz  $\mathbf{d}_M(\mathbf{x}_A, \mathbf{x}_B) = \sum_{j=1}^p |\mathbf{x}_{Bj} - \mathbf{x}_{Aj}|$

Minkowski Distanz  $\mathbf{d}_{Mink}(\mathbf{x}_A, \mathbf{x}_B) = \left( \sum_{j=1}^p (\mathbf{x}_{Bj} - \mathbf{x}_{Aj})^m \right)^{\frac{1}{m}}$

Kosinus des Winkels  $\mathbf{COS} \alpha = \frac{\mathbf{x}_A^T \mathbf{x}_B}{\sqrt{(\mathbf{x}_A^T \mathbf{x}_A)(\mathbf{x}_B^T \mathbf{x}_B)}} = \frac{\mathbf{x}_A^T \mathbf{x}_B}{\|\mathbf{x}_A\| \cdot \|\mathbf{x}_B\|}$



Mahalanobis Distanz  $\mathbf{d}_{Mahal}(\mathbf{x}_A, \mathbf{x}_B) = \left[ (\mathbf{x}_B - \mathbf{x}_A)^T \Sigma^{-1} (\mathbf{x}_B - \mathbf{x}_A) \right]^{\frac{1}{2}}$   $\Sigma^{-1}$  inverse Kovarianzmatrix

$\mathbf{d}_{Mahal}(\mathbf{x}_i, \boldsymbol{\mu}) = \left[ (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right]^{\frac{1}{2}}$  für  $i = 1$  bis  $n$ ,  $\boldsymbol{\mu}$  ... Zentrum der Verteilung

# Linearkombinationen

$$u = b_1 x_1 + b_2 x_2 + \dots + b_p x_p \quad u = x^T b$$

## Hauptkomponenten (PC)

$U = XB$   $X$  ... projizierte Datenwerte  $\rightarrow (n \times p)$   $B$  ... (Spalten = Ladungen der PC  $\rightarrow (p \times p)$ )

$$b_j = (b_{1j}, \dots, b_{pj})^T \rightarrow \text{Max}(\text{Var}(u_j)) \text{ \& } b_j^T b_l = 0 \text{ \& } b_j^T b_j = 0 \quad \text{shape}(b_j) = p \times 1$$

$$\text{Var}(u_j) = \text{Var}(x_1 b_{1j} + \dots + x_p b_{pj}) = b_j^T \text{Cov}(x_1, \dots, x_p) b_j = b_j^T \Sigma b_j$$

$\Sigma$  ... *theo. Kovarianzmatrix*

Lagrange Problem  $\phi_j = b_j^T \Sigma b_j - \lambda_j (b_j^T b_j - 1)$   $j = 1$  bis  $p$   $\lambda_j$  ... *Langrange Parameter*

$$\text{Var}(u_j) = b_j^T \Sigma b_j = b_j^T \lambda_j b_j = \lambda_j b_j^T b_j = \lambda_j$$

Anteil der erste  $k$  PCs  $\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p}$

# Clusteranalyse

k-means Zentroide  $\bar{x}_j = \frac{1}{n_j} \sum_{i \in I_j} x_i$   $j = 1$  bis  $k$  Zielfunktion:  $\sum_{j=1}^k n_j \sum_{i \in I_j} \|x_i - \bar{x}_j\|^2 \rightarrow \min$

## Hierarchische Clustermethoden

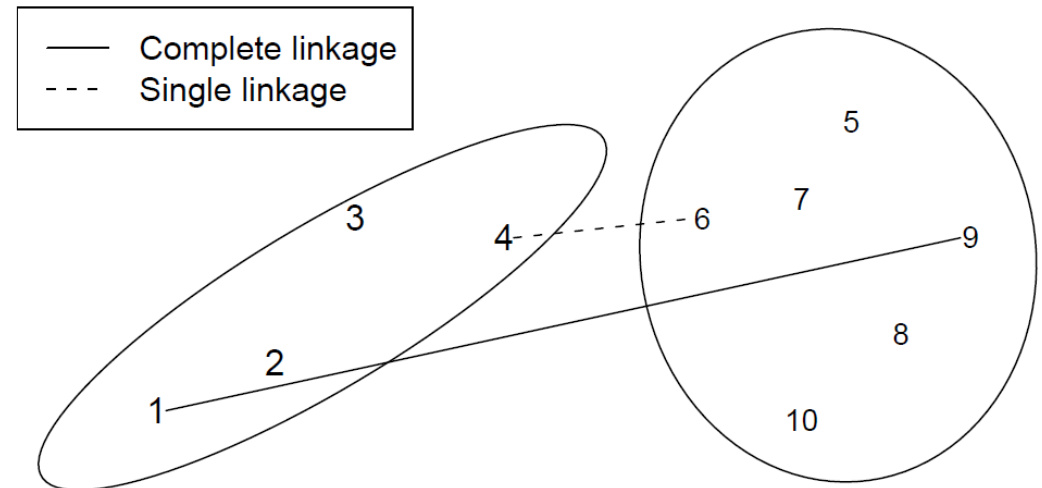
Complete linkage:  $\max_{i \in I_j, i' \in I_j} d(x_i, x_{i'})$

Single linkage:  $\min_{i \in I_j, i' \in I_j} d(x_i, x_{i'})$

Average linkage:  $\text{average } d(x_i, x_{i'})_{i \in I_j, i' \in I_j}$

Centroid Methode:  $d(\bar{x}_j - \bar{x}_{j'})$

Ward Methode:  $d(\bar{x}_j - \bar{x}_{j'}) \frac{\sqrt{2n_j n_{j'}}}{\sqrt{n_j + n_{j'}}}$



## Fuzzy Clustering

c-means Zentroide  $\bar{x}_j = \frac{\sum_{i=1}^{n_j} u_{ij}^2 x_i}{\sum_{i=1}^{n_j} u_{ij}^2}$  Zielfunktion:  $\sum_{j=1}^k \sum_{i \in I_j} u_{ij}^2 \|x_i - \bar{x}_j\|^2 \rightarrow \min$   $u_{ij}$  Zugehörigkeit

# Gütemaße

Within-Cluster Sum-of-Squares (für Varianz)

$$W_k = \sum_{j=1}^k \sum_{i \in I_j} \|\bar{x}_j - \bar{x}\|^2$$

Between-Cluster Sum-of-Squares (für Heterogenität)

$$B_k \sum_{j=1}^k \|\bar{x}_j - \bar{x}\|^2 \text{ mit } \bar{x} = \frac{1}{k} \sum_{j=1}^k \bar{x}_j$$

Calinski-Harabasz-Index

$$CH_k = \frac{B_k / (k - 1)}{W_k (n - k)}$$

Hartigan-Index

$$H_k = \log \left( \frac{B_k}{W_k} \right)$$



# Diskriminanzanalyse

$$\phi_j(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma_j)}} \exp \left\{ -\frac{(\mathbf{x} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x} - \mu_j)}{2} \right\}$$

p ... Features,  $\Sigma$  ... Kovarianz, j ... Gruppe,

$\phi_j$  .. Dichtefunktion pro Gruppe

$$P(G = j | \mathbf{x}) = \frac{\phi_j(\mathbf{x}) p_j}{\sum_{l=1}^k \phi_l(\mathbf{x}) p_l}$$

$p_j$  ... a-priori Wahrscheinlichkeit pro Gruppe. G ... Gruppenzugehörigkeit

x wird j zugeteilt wenn gilt:  $\log \left( \frac{P(G = j | \mathbf{x})}{P(G = l | \mathbf{x})} \right) = \log \left( \frac{\phi_j(\mathbf{x}) p_j}{\phi_l(\mathbf{x}) p_l} \right) = \log \left( \frac{\phi_j(\mathbf{x})}{\phi_l(\mathbf{x})} \right) + \log \left( \frac{p_j}{p_l} \right) > 0$

## Lineare Diskriminanzanalyse (LDA)

es gilt:  $\Sigma_1 = \dots = \Sigma_k = \Sigma$

lineare Diskriminanzfunktion:

$$\delta_j(\mathbf{x}) = \mathbf{x}^T \Sigma^{-1} \mu_j - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j + \log p_j$$

## Quadratische Diskriminanzanalyse (QDA)

$$\delta_j^{(q)}(\mathbf{x}) = -\frac{1}{2} \log \left( \det(\Sigma_j) \right) - \frac{1}{2} (\mathbf{x} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x} - \mu_j) + \log(p_j)$$