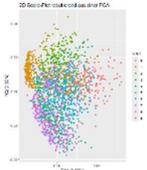
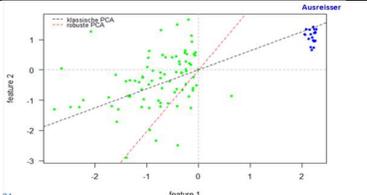
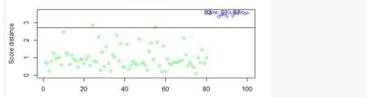
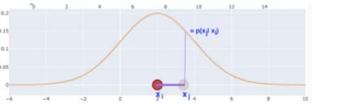


<p>Data Mining</p> <p>Anwendung geeigneter Methoden zur Entdeckung von Strukturen & Zusammenhängen in Datenmengen</p> <ul style="list-style-type: none"> - Für grosse Datenbestände - Anz. Merkmale kann Anz. Beobachtungseinh. übersteigen - Charakteristisch, Daten wurden nicht zu einem bestimmten Zweck erhoben auch anwendbar, wenn traditionelle Techniken versagen - Garbage in, garbage out (GIGO) 	<p>Data Preparation: Umgang mit fehlenden Werten / Ausr., Standardisierung, Transformation, aggregieren und selektieren von Var.</p> <p>Modeling: Auswahl Data Mining Verfahren, Prognosen, Auswertungen</p> <p>Evaluation: Statistische Evaluierung des Modells, qualitative Beurteilung</p>
--	--

<p>Allgemein: Multivariate Daten – p Dimensionen</p> <p>Auch bei überwachten Verfahren häufig hilfreich</p> <ul style="list-style-type: none"> - Pre-processing Schritte vor Reduktion - Kann helfen, irrelevante Variablen eliminieren oder Noise in den Daten reduzieren - Zeit und Rechenpower reduzieren - Gewisse Klassifikationsmodelle können nicht mit hochdimensionalen Daten umgehen (Regression) <p>Modus: Multivariater Modus kann nicht über Moduswerte der Randverteilung abgeleitet werden</p> <p>Median: Keine allgem. Definition für multiv. Median, am weitesten verbreitet ist L1-Median (Punkt mit minim. Summe der absoluten Abstände zu allen anderen Punkten), komponentenweiser Median ergibt nicht dasselbe R-Code: L1median (package: robustX)</p>	<p>Mittelwert: Komponentenweiser Mittelwert entspricht dem multivariaten Schwerpunkt</p> <p>Ausreisser: Multivariater Ausreisser (grün) ist nicht zwingend auch ein Ausreisser in der Randverteilungen (blau)</p> <p>→ Hochdimensionale Räume haben häufig extreme Merkmalswerte → effektive Ausreisser analysieren → Alle Punkte haben fast denselben Abstand zueinander (mit dist(dat) Distanzmatrix berechnen)</p> <p>Benötigte Beobachtungen:</p> <table border="1"> <tr> <th>Dimension</th> <th>Anzahl Beobachtungen</th> </tr> <tr> <td>1</td> <td>10</td> </tr> <tr> <td>2</td> <td>100</td> </tr> <tr> <td>p</td> <td>10^p</td> </tr> </table>	Dimension	Anzahl Beobachtungen	1	10	2	100	p	10^p	<p>Datenübersicht</p> <pre>str(dat) dim(dat) table(dat\$zielvar) → Übersicht Verteilung table(dat\$zielvar)/sum(dat\$zielvar) → Anteil</pre> <p>Visualisierung:</p> <pre>boxplot(dat) boxplot(var ~ zielvar, data = dat) pairs(dat[,2:4], col = c("red", "blue"))[dat\$zielvar] → univariate Beziehung</pre> <pre>moaicplot(~ var + zielvar, data = train_data) barplot(table(dat\$zielvar), las = 2)</pre> <p>Skalieren, wenn Einheiten nicht gleich sind, wenn alle Variablen gleiches Gewicht haben sollen Zentrierung, wenn bereits alle Einheiten & Varianzen gleich / ähnlich sind, um Ausreisser eher hervorzuheben</p>
Dimension	Anzahl Beobachtungen									
1	10									
2	100									
p	10^p									

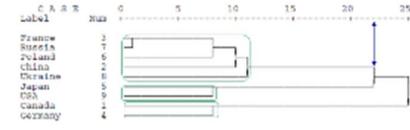
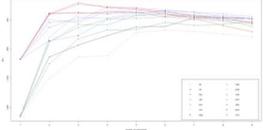
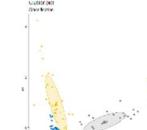
<p>Ähnlichkeits- und Distanzmasse</p> <p>Ähnlichkeitsmass s (similarity)</p> <p>→ quantifiziert die Übereinstimmung → je grosser der Wert desto grösser die Ähnlichkeit</p> <p>Distanzmasse für quantitative Variablen</p> <p>Axiome Distanzmass</p> <p>Eigenschaften (reellwertige Distanzfunktion):</p> <ol style="list-style-type: none"> 1. $d(O_1, O_2) \geq 0$ → nicht negativ 2. $d(O_1, O_1) = 0$ → Dist. zu sich selbst = 0 3. $d(O_1, O_2) = d(O_2, O_1)$ → symmetrisch <p>Metrische Distanzen erfüllen zusätzlich folgende Eigenschaften:</p> <ol style="list-style-type: none"> 4. $d(O_1, O_2) = 0$, genau dann wenn $O_1 = O_2$ 5. $d(O_1, O_2) \leq d(O_1, O_3) + d(O_3, O_2)$ → Dreiecksungleichung <p>Minkowski Distanz</p> <p>→ $d_{ij}(o_i, o_j) = (\sum_{k=1}^p o_{ik} - o_{jk} ^r)^{\frac{1}{r}}$</p> <p>Euklidische Distanz (L2-Metrik)</p> <p># dist() auf Distanzmatrix anwenden → nur bei num Variablen sonst as.matrix(daisy(data, type=list(ordratio)))[c("A", "C"), "B"] # Parameter: ordratio behandelt Variablen mit Verhältnisskala ordinal, type=list(ord...) zwingend davor → Distanz mit Pythagoras berechnen, direkte Distanz Mass vom Range der Variablen abhängig → bei Bedarf für Rob. Skal. im hochdimensionalen Raum: Entspricht PCA</p>	<p>Manhattan Distanz (L1-Metrik)</p> <p>→ $d_{ij}(o_i, o_j) = \sum_{k=1}^p o_{ik} - o_{jk}$</p> <p>Für hochdimensionale Daten besser geeignet als euklidische Distanz da auf grosse Abstände nicht quadratisch eingegangen wird, nicht auf direktem Weg «um Blocks»</p> <p>Distanzmasse für binäre Daten</p> <p>Die Objekte O_1 und O_2 haben nur binäre Attribute. Z. B. Geschlecht (f/m), Fahrlizenz (yes/no), Nobelpreisträger (yes/no), ...</p> <p>Beispiel:</p> <table border="1"> <tr> <td>$O_1 = 1000000000$</td> <td>Objekt o_1</td> <td>0</td> <td>1</td> <td></td> </tr> <tr> <td>$O_2 = 0000001001$</td> <td>Objekt o_2</td> <td>$M_{00} = 7$</td> <td>$M_{01} = 1$</td> <td>$M_{10} = 8$</td> </tr> <tr> <td></td> <td></td> <td>$M_{10} = 2$</td> <td>$M_{11} = 0$</td> <td>$M_{11} = 2$</td> </tr> <tr> <td></td> <td></td> <td>$M_{10} = 9$</td> <td>$M_{11} = 1$</td> <td>$M_{11} = 10$</td> </tr> </table> <p>Symmetrische Variablen: Beide Levels ähnlich häufig (z.B. Geschlecht):</p> <ul style="list-style-type: none"> - Simple-Matching: $s(o_1, o_2) = \frac{M_{11} + M_{00}}{(M_{01} + M_{10} + M_{11} + M_{00})}$ - Canberra-Metrik: $d(o_1, o_2) = 1 - s(o_1, o_2)$ <p>Asymmetrische Variablen: Levels sehr unterschiedlich</p> <ul style="list-style-type: none"> - Jaccard Koeffizient $s(o_1, o_2) = \frac{M_{11}}{(M_{01} + M_{10} + M_{11})}$ - Ohne Variabel mit der höchsten Anzahl - Jaccard-Metrik $d(o_1, o_2) = 1 - s(o_1, o_2)$ 	$O_1 = 1000000000$	Objekt o_1	0	1		$O_2 = 0000001001$	Objekt o_2	$M_{00} = 7$	$M_{01} = 1$	$M_{10} = 8$			$M_{10} = 2$	$M_{11} = 0$	$M_{11} = 2$			$M_{10} = 9$	$M_{11} = 1$	$M_{11} = 10$	<p>Mehr als 2 Levels (nominale und ordinale Daten)</p> <p>Objekte O_1 und O_2 haben p nominale Merkmale:</p> <p>Sneath-Übereinstimmungskoeffizient</p> <p>$s(o_1, o_2) = \frac{\# \text{Übereinstimmungen}}{p}$</p> <p>$d(o_1, o_2) = 1 - s(o_1, o_2)$ → nicht Übereinstimmungen</p> <p>Falls geordnete Kategorien mit K Stufen (ordinale Daten)</p> <p>Kategorien die Werte 1 bis K zu ordnen → x_0 Transformieren: $z_0 = \frac{x_0 - 1}{K - 1} \in [0, 1]$ → Distanzen für quantitative Daten</p> <p>Mischung von Variablentypen - Distanzmass von Gower</p> <p>Distanzmasse zwischen 0 und 1 werden für alle Variablen verwendet.</p> <ul style="list-style-type: none"> - K-tes Merkmal ist numerisch: Rk: Range der Variable k für alle Objekte $d(o_i, o_j) = \frac{ o_{ik} - o_{jk} }{Rk}$ - K-tes Merkmal ist binär, nominal: Canberra-Metrik bzw. Jaccard-Metrik für asymmetrische Variablen - K-tes Merkmal ist ordinal: Transformation in Ränge und diese anschliessend wie numerische Grösse behandeln <p>Aggregation über alle Variablen</p> <p>$d_{ij} = \frac{1}{p} \sum_{k=1}^p d_{ij}^k$ → Grösse zwischen 0 und 1</p>
$O_1 = 1000000000$	Objekt o_1	0	1																			
$O_2 = 0000001001$	Objekt o_2	$M_{00} = 7$	$M_{01} = 1$	$M_{10} = 8$																		
		$M_{10} = 2$	$M_{11} = 0$	$M_{11} = 2$																		
		$M_{10} = 9$	$M_{11} = 1$	$M_{11} = 10$																		

Unüberwachtes Lernen		Dimensionsreduktion → Zweck: Visualisierung hochdimensionaler Daten				
		- Strukturen erkennen - Ausreisser ermitteln				
Methode	Ziel	Voraussetzungen	Vorgehen	Plot	R-Code	
PCA (Hauptkomponentenanalyse)	Achsen finden, welche multidim. Daten erklären - Variablen sollen unkorreliert sein - Grossen Teil der Totalvarianz abdecken	<ul style="list-style-type: none"> • Datenmatrix • Originalvariablen stark korreliert • kat. Variablen als factor() • Numerische Variablen • Annahme: Daten normalverteilt und frei von Ausreissern 	<ol style="list-style-type: none"> 1. Koordinatensystem in Schwerpunkt der Daten verschieben 2. Rotation bis Varianz entlang der ersten Hauptkomponente am grössten ist → Grösste Teil der restlichen Varianz entlang der zweiten Hauptkomponente (steht orthogonal zur Ersten) <p>Rotation: Ergebnis durch Multiplikation der zentrierten Datenmatrix X mit Rotationsmatrix A, → Matrizen</p>	<p>biplot() → rote Pfeile Plot, autoplot() → 2D-Score-Plot</p> <p>Eher auf Anteil verlassen</p> <p>PCA mit 9 Variablen</p> <p>Score-Plot für PCA mit 9 Variablen</p> <p>Score-Plot: Ellbogenmethode</p>	<pre>prcompmt(data, scale = TRUE)</pre> <p>library(ggfortify), autoplot(...) → Scoreplot, biplot(...)</p> <p>Überprüfung Streuung: boxplot(ccf, 2:x) # x: (alle Var - Zielvar)</p> <p>Falls unterschiedliche Streuung der Variablen: pca_klassisch <- prcomp(creditcard[, 2:30], scale=TRUE) plot(pca_klassisch\$x[, 1], pca_klassisch\$x[, 2], # Plot Anteil (siehe links) col=c("red", "blue"))[as.factor(creditcard\$Class)], pch=20, las=1, xlab="PC1", ylab="PC2"] var <- pca_klassisch\$sdev^2 # 80% der Var ermitteln </p>	

<p>Robuste PCA</p>	<p>PCA Methode direkt mit robusten Schätzern durchführen</p>	<p>PCA Vorteile: - Nutzen besonders hoch, wenn Originalvariablen stark korreliert sind - Wichtigsten Hauptkomponenten können für weiterführende Analysen verwendet werden</p> <p>PCA Nachteile: - Abhängig von Skalierung - Hauptkomponenten schwierig zu identifizieren - Anzahl relevante Hauptkomponenten nicht eindeutig</p> <p>NUR FÜR NUMERISCHE VARIABLEN</p>	<p>können nur multipliziert werden, wenn ihre Dimensionen kompatibel sind (siehe LA Regeln)</p> <p>Varianz der j-Variabel: In Diagonalelementen enthalten, die anderen Werte entsprechen der Kovarianz der j-ten und k-ten Variabel</p> <p>Kovarianzmatrix: Quadratisch & symmetrisch</p> <p>Falls Variablen korrelieren: Werte ausserhalb der Diagonale haben Wert ≠ 0, falls keine Korrelation Werte = 0</p> <p>Autoplot: library(ggfortify) x\$label <- as.factor(x\$label) res_pca <- prcomp(x[, -1]) autoplot(res_pca, data = x, colour = 'label', label.size = 3) + ggtitle("2D Score-Plot resultierend aus einer PCA")</p>  <p>res.pca <- prcomp(diabetes[, -1]) plot(res.pca\$x, pch=20, main="Klassen", col = diabetes\$class) legend("topright", legend = levels(diabetes\$class), col = 1:3, pch=20)</p>	<p>Dimensionsreduktion Verwenden nur die ersten paar Hauptkomponenten → Daten sollen durch weniger Komponenten möglichst gut beschrieben werden Gütekriterium: ~ 80 % der totalen Varianz</p>  <p>Es gibt zwei Arten von Ausreissern. • Ausreisser im PCA-Raum • Score-Distanz $SD_i = \sqrt{\sum_{j=1}^k x_{ij}^2}$ (optionale Cutoff-Linie $\sqrt{x_{1,0.975}}$)</p>  <p>→ Oberhalb Linie = potentielle Ausreisser</p>	<pre>var_cum <- cumsum(var)/sum(var) plot(1:length(var),var_cum, ylim=c(0,1)) # Ergebnis 80 % Varianz plotten pc_number <- min(which(var_cum > 0.8)) abline(h=0.8,y=pc_number) grid() # Gitternetzlinien min(which(var_cum > 0.8)) # Dim.Reduktion: Num-Ausgabe</pre> <pre>library("MASS") & library("rccov"), library(robustbase) pca.rob <- PcaHubert(creditcard[, 2:30], scale = TRUE, k = 2) # Argument k: Entspricht der Anzahl berechneten Hauptkomponenten. u dieser Dimension wird auch die orthogonale Distanz berechnet. plot(pca.rob@scores[, 1], pca.rob@scores[, 2], col = c("red", "blue")[as.factor(creditcard\$class)], pch=20, las=1, xlab="PC1", ylab="PC2")</pre> <p>Orthog. Ausreisser (Punkte mit grossem Rekonstruktions-Fehler) in Plot mit orth. Distanzen prüfen & Cutofflinie einzeichnen: n <- nrow(creditcard) plot(x=1:n,y=pca.rob@od,pch=16,col=c(rgb(1,0,0,0.6),"blue")) [as.factor(creditcard\$class)] points(x= which(c_card\$class==1), y=p_rob@od[c_card\$class==1],col='blue',pch=16) # points stellt sicher, dass alle Betrugsfälle erkennbar sind (nicht überdeckt) abline(h=pca.rob@cutoff.od) # Zeilen, die nach diesem Kriterium Ausreisser sind: outlier.rows <- which(pca.rob@od >= pca.rob@cutoff.od) table(ifelse((pca.rob@od >= pca.rob@cutoff.od) == TRUE, "Ausreisser", "Kein Ausreisser"), ifelse(cc\$class == 1, "Betrug", "Kein Betrug")) # Konf.matrix Abdeckung in % der Hauptkomponenten: summary(pca)</p>
<p>MDS</p>	<p>Beziehung zwischen Obj. basierend Ähnlichkeiten oder Distanzen im niedr.dim. Raum visualisieren. (Gleiches wie PCA, aber mit Distanzen) ► metrische MDS ► nicht metrische oder ordinale MDS</p>	<ul style="list-style-type: none"> • Distanzmatrix • Wenn kat. Variablen, dann dasyj() anwenden 	<p>PCA und metrische MDS sind äquivalent, wenn euklidische Distanzen verwendet werden (metrische Distanzen, mittels Eigenwertzerlegung werden Koordinaten für Visualisierung ermittelt). MDS ist (meistens) auch für nicht euklidische Distanzen geeignet. Metrische MDS benötigen nur Distanzen.</p> <p>Falls die Daten nicht metrisch sind: Ordinale MDS durchführen</p>	<p>normaler plot(): res.cmd <- cmdscale(eurodist, k=2) #Metrisches MDS plot(res.cmd, pch="m") text(res.cmd, labels = rownames(res.cmd))</p> <p>ggplot() mit Einfärbung: library(ggplot2), library(plotly) res_mds <- cmdscale(dist_matrix23) df <- data.frame(res_mds) plot_data <- as.data.frame(df) plot_data <- cbind(plot_data, Meta23)</p> <p>P<-ggplot(plot_data, aes(x=X1, y=X2,label=Name,colour=Fraktion))+ geom_text(size=3)+ geom_point() + scale_color_manual(values=c("green","yellow","orange"))</p> <p>ggplotly(P) # Möglichkeit zur Identifizierung einzelner Leute</p>	<p>Ordinale nicht-metrische MDS: isoMDS() funktioniert nicht mit Distanzen = 0 Entfernen dieser Distanzen (je nach Format der Daten): dist_b<-as.dist(as.matrix(dist_m)[-which(colnames(as.matrix(dist_m))=="Martin"), -which(colnames(as.matrix(dist_m))=="Martin")]) #Spalte Meta_b <- Meta23[-which(Meta23\$Name == "Martin"),]</p> <p>Ordinale oder nicht-metrische MDS: library(MASS) res.iso <- isoMDS(eurodist, k=2) plot(res.iso\$points, pch="m") text(res.iso\$points, labels = rownames(res.iso\$points))</p>
<p>t-SNE</p>	<p>Versucht Nachbarschaftsbeziehung n im hoch-dimensionalen Raum zu erhalten und in niedriger Dimension darzustellen.</p>	<ul style="list-style-type: none"> • Datenmatrix: num • Distanzmatrix: num./kat. <ul style="list-style-type: none"> • Wenn nötig, PCA durchführen • richtige Wahl von Parameter durch Ausprobieren <p>→ Ist t-SNE / UMAP besser als PCA: Daten liegen evtl. nicht in einer Hyperebene → Struktur komplexer</p>	<p>Berechnung Ähnlichkeit von einem Punkt zu allen anderen Punkten (euklidische Distanz), Punkte werden danach in eine Normalverteilung X-Achse eingetragen (sofern σ^2 gegeben ist) → Bedingte W'keit Je näher x_j an x_i (roter Punkt) ist, desto ähnlicher sind sich die Punkte → Vorgang für alle Punkte wiederholen</p>  <p>Mass für erw. Überraschung: $H(P_i) = \sum_j p_{ji} \log_2 \left(\frac{1}{p_{ji}} \right)$</p> <p>Crowding-Problem: Dimensionsreduktion führt zu Struktur Verfälschung -> Durch t-Verteilung gelöst Wird die bed. W'keit zwischen 2 Punkten berechnet (xi zu xj & xj zu xi) werden Werte unterschiedlich sein: Sie stammen aus zwei Verteilungen</p>	<p>Vorsicht: Resultate von t-SNE können auch irreführend sein, Clustergrößen & Abstände haben keine Bedeutung → lediglich Darstellung von lokalen Nachbarschaften, mehrfaches Ausführen kann zu lokalen Optimas führen, kann nicht auf neue Punkte angewandt werden</p> <p>Vergleich Plot PCA & t-SNE: plot_data = data.frame(res_tsne\$Y) colnames(plot_data) <- c("x1", "x2")</p> <p>par(mfrow=c(1,2)) plot(res_tsne\$Y[, 1], res_tsne\$Y[, 2], type="n", main = "t-SNE", xlab = "t-SNE1", ylab = "t-SNE2") text(res_tsne\$Y[, 1], res_tsne\$Y[, 2], labels=row.names(abst), cex=0.8)</p> <p>res_pca <- prcomp(abst, scale=FALSE) plot(res_pca\$x[, 1], res_pca\$x[, 2], type="n", main = "PCA", xlab = "PC1", ylab = "PC2") text(res_pca\$x[, 1], res_pca\$x[, 2], labels=row.names(abst), cex=0.8)</p>	<p>library(Rtsne) Empfehlung: 1. PCA 2. t-SNE auf die ersten (50) Hauptkomponenten</p> <p>perplexity: 3 - perplexity < nrow (X) - 1 → Mit steigender Perplexity nimmt Varianz zu Perplexity ~ Anzahl lokale Nachbarn (Werte zwischen 5 – 30 sinnvoll)</p> <p>iterationen: ausreichend gross wählen Kategorielle Variablen: z.B. Label = as.factor(x\$label) Label für Rtsne / UMAP nicht verwenden, in Plot jedoch wieder hinzufügen.</p> <pre>res_tsne <- Rtsne(X=dist_matrix23[, -1], dims = 2, is_distance=TRUE, pca = FALSE, perplexity = 10, max_iter = 1000) #Zielvar. entfernen</pre> <p>plot_data <- as.data.frame(res_tsne\$Y) plot_data <- cbind(plot_data, Meta23) plot_data\$label <- dist_matrix23\$label</p> <p>ggplot(plot_data, aes(x = V1, y = V2, label= Name, col=label)) + geom_text(size=3) + geom_point() + scale_color_manual(values=c("green", "yellow", "orange", "blue"))</p>

<p>UMAP</p>	<ul style="list-style-type: none"> • Dimensionsreduktion für mehr als 2-3 Dimensionen • Schneller als t-SNE • Oftmals globale Strukturen besser erhalten als t-SNE • Neue Daten können zur Projektion hinzugefügt werden (wenn Datenmatrix vorhanden) • Kann auf Rohdaten angewandt werden (keine vorgängige PCA) 	<ul style="list-style-type: none"> • Datenmatrix: num • Distanzmatrix: num./kat. • numerische Werte kann auf Rohdaten ohne PCA ausgeführt werden. • richtige Wahl von Parameter durch Ausprobieren 	<ul style="list-style-type: none"> • Abstand wird über Exponentialverteilung bestimmt, Einzeichnung auf der exp. Abfallenden Kurve -> grössere Abst. der Cluster • Mittelung der Ähnlichkeiten in beide Richtungen • Summe der unskal. Ähnlichkeiten (Dichte): $\log(\text{Anz. Nachbarn})$ • Minimierung der Kostenfunktion (binäre Kreuzentropie) 	<p>Clustergrößen haben keine Bedeutung (entscheidend: lokale Abstände), Abstände zwischen Clustern sind nicht immer sinnvoll (können müssen aber keine Bedeutung haben), zufälliges Rauschen sieht nicht immer zufällig aus, Cluster oft besser getrennt als in t-SNE</p>	<pre>library(umap) res_umap <- umap(d=as.matrix(dist_matrix23), input = "dist", n_neighbors=25, min_dist=0.2) #Weitere Parameter: random_state umap.defaults plot_data <- as.data.frame(res_umap\$layout) colnames(plot_data) <- c("x1", "x2") plot_data <- cbind(plot_data, Meta23) ggplot(plot_data, aes(x=x1, y=x2, label=Name, colour=Fraktion)) + geom_text(size=3) + geom_point() + scale_color_manual(values=c("green", "yellow", "orange")) + ggtitle("2D UMAP Visualisierung")</pre>
--------------------	--	--	---	---	--

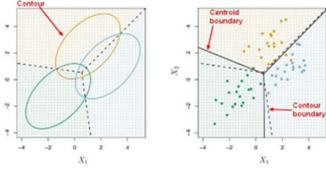
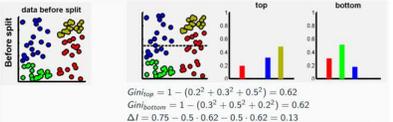
<p>Clustering</p>		<ul style="list-style-type: none"> • Einteilung von Objekten in homogene Gruppen / Cluster • Innerhalb der Cluster: Objekte sind untereinander ähnlicher als zu Objekten aus anderen Clustern <p>Vorgehen:</p> <ol style="list-style-type: none"> 1. Explorative Datenanalyse / Datenbereinigung 2. Skalierung der Daten 3. Wahl eines Cluster-Algorithmus 4. Evaluierung Anzahl Cluster <p>Visualisierung: Versch. Farben in 2D (PCA, MDS, t-SNE, UMAP)</p>			
<p>Methode</p>	<p>Ziel</p>	<p>Voraussetzung</p>	<p>Vorgehen / Anwendung</p>	<p>Plot</p>	<p>R-Code</p>
<p>K-Means</p>	<p>Unterteilung der Beobachtungen in k homogene Cluster – totale Variation (WCV) innerhalb der Cluster über alle k Cluster so klein wie möglich ist.</p> $WCV(C_k) = \frac{1}{\#C_k} \sum_{i,j \in C_k} \sum_{j=1}^p (x_{ij} - x_{rj})^2$ <p>WCV = eukl. Dist. $\hat{2}$ zw. Objekten $\#C_k$ = Anz. Beob. im Cluster k p = Anz. Features (Dimensionen)</p> <p>→ Liefert lokale Lösung, Abhängig von Startkonfiguration → mehrmals durchführen bis kleinstes WCV gefunden</p> <p>Nicht anwenden bei:</p> <ul style="list-style-type: none"> - Verrauschten Daten / Ausreißern - Schwierig bei hochdimensionalen Daten (nicht aussagekräftig) 	<ul style="list-style-type: none"> • Datenmatrix • kat. Variable als factor <p>• Anz. Cluster muss vorgegeben werden</p> <ul style="list-style-type: none"> • Festlegung der Anzahl Cluster durch: <ul style="list-style-type: none"> - Elbogen-Methode - Silhouetten-Methode - Domain-Knowledge 	<ol style="list-style-type: none"> 1. Beobachtungen zufällig einer Gruppe zuordnen 2. Clusterzentrum bestimmen 3. Alle Beobachtungen werden dem Clusterzentrum zugeordnet zudem die kl. Distanz <p>Wiederholen 2. & 3. bis fertig</p> <p>Qualitätskontrolle: Silhouetten <i>Silhouetten – Koeffizient, $sil_i = \frac{b_i - a_i}{\max(a_i, b_i)}$</i> a = mittl. Distanz zwischen Punkt i und allen anderen Punkten aus dem gleichen Cluster b = mittl. Distanz zwischen i und allen Punkten des nächsten Nachbarclusters</p> <p>• $bi \gg ai \rightarrow$ ben. Cluster weit weg sil. $\rightarrow 1$ • $bi \approx ai \rightarrow$ ben. Cluster nahe sil. $\rightarrow 0$ • Beobachtungen im Mittel näher beim benachbarten Cluster $\rightarrow < 0$</p> <p>Faustregel für mittlere Breite in einem Cluster: 0.70 < sil < 1.00: Gute Struktur gefunden 0.50 < sil < 0.70: vernünftige Strukturen 0.25 < sil < 0.50: Schw. Strukt. \rightarrow w. Unters. -1 < sil < 0.25: Forget it!</p>	<p>R-Code Silhouetten Plot: <code>library(cluster)</code> <code>cl <- kmeans(x, centers=3)</code> <code>plot(silhouette(x=cl\$cluster, dist=dist(x)))</code> <i># Bei PAM: x = res_pam\$clustering</i></p> <p>Anzahl Cluster mit mittl. Silhouettenbreite finden: auch für pam <code>sil <- rep(0, 15)</code> for (i in 2:15){ <code>sil[i] <- summary(silhouette(x=kmeans(x, k=i)\$cluster, dist=dist(x)))\$avg.width</code> <i># Bei PAM: x = pam(dat,</i> <code>k=i)\$cluster</code> <code>plot(x=1:15, sil, type = "b", las=1, xlab = "Anzahl Cluster", ylab = "mittlere Silhouetten-Breite")</code></p> <p>Opt. Clusteranz. mit Within-cluster Variation (Elbogen-Methode) <code>wcv <- rep(0, 8)</code> # Initialisierung for (i in 1:8) <code>wcv[i] <- sum(kmeans(abst, centers = i)\$withinss)</code> <code>plot(1:8, wcv, type = "b", xlab = "Anzahl Cluster", ylab = "within-cluster Variation", las=1)</code></p> <p>Alternative: <code>library(factoextra)</code> anstelle von kmeans kann auch pam eingesetzt werden -> braucht aber Distanzmatrix <code>fviz_nbclust(x, kmeans, method="silhouette")</code> # <i>method = "wss"</i></p> <p>nur für pam: <code>library(fpc)</code> \rightarrow Benötigt Distanzmatrix <code>cl_pamk <- pamk(x, krange=2:5)</code> <code>cl_pamk\$nc</code> # <i>Output: Anz. Cluster</i></p>	<pre>x <- rbind(matrix(rnorm(100, sd = 0.3), ncol = 2) # Samples generieren matrix(rnorm(100, mean = 1, sd = 0.3), ncol = 2) matrix(rnorm(100, mean = 1.6, sd = 0.3), ncol = 2)) colnames(x) <- c("x", "y") cl <- kmeans(x, centers=3) # K-Means anwenden plot(x, col = cl\$cluster, las=1) # Plot points(cl\$centers, bg = 1:3, col = "purple", ch = 22, cex = 2) cl\$withinss # Abstände zwischen den Clustern cl\$tot.withinss # Total Abstand Bei > 2 Dimensionen: Vis. mit 2D-Plot mit PCA, MDS, t-SNE oder UMAP Anzahl Cluster (k, param: center in kmeans) finden: wcv <- rep(0, 6) # Range für wcv generieren for (i in 1:6) wcv[i] <- sum(kmeans(x, centers = i)\$withinss) plot(1:6, wcv, type = "b", xlab = "Anzahl Cluster", ylab = "within-cluster Variation", las=1) # Plot, Anzahl Cluster ablesen: k nach dem letzten grossen Abfall im Plot Alternative: library(factoextra) fviz_nbclust(x, kmeans, method="wss", k.max = 10)</pre>
<p>K-Medoids / PAM</p>	<p>Repräsentative Objekte finden (Medoids) Ausgehend von einer Distanzmatrix wird über alle Objekte die Gesamtdistanz der Objekte zu den Objekten im gleichen Cluster minimiert</p> <ul style="list-style-type: none"> - Langsam für grosse Daten \rightarrow <code>clara()</code> - Anzahl Cluster muss vorgegeben sein 	<ul style="list-style-type: none"> • Distanzmatrix <p>• Festlegung der Anzahl Cluster durch: - Silhouetten-Methode - Domain-Knowledge</p> <p>+ Robuster gegenüber Ausreißern + Kann mit allen Distanzmassen umgehen</p>	<ol style="list-style-type: none"> 1. k Objekte zufällig als Clusterzentrum definieren 2. Jedes Objekt dem nächsten Clusterzentrum zuweisen 3. Kostenfunktion wird evaluiert 4. Iterativ wird für jeden Cluster untersucht, ob das vertauschen des Medoids mit einem nicht Medoid-Objekt zu einem verbesserten Kostenwert führt. Falls zumindest ein Medoid vertauscht wurde, beginnt der Algorithmus wieder bei Schritt 2. Ansonsten endet der Algorithmus. 	<p>Silhouetten Plot mit Balken: <code>res.km2 <- kmeans(abst, 2)</code> <code>sil_cl2 <- silhouette(x=res.km2\$cluster, dist=dist(abst))</code> <code>rownames(sil_cl2) <- rownames(abst)</code> <code>plot(sil_cl2, cex.names=0.6)</code></p>	<p>R-Code: <code>library(cluster)</code> <code>cl_pam <- pam(x, k=3)</code> <code>plot(x, col = cl_pam\$clustering) # col auch so bei UMAP, tSNE MDS</code> <code>points(cl_pam\$medoids, col = 1:3, pch = 8, cex = 2)</code> <code>cl_pam\$objective</code></p> <p>Bestimmung Anzahl Cluster über mittl. Silhouetten-Breite</p> <p>Repräsentative Beobachtungen: <code>res_pam\$medoids</code></p>
<p>Hierarch. Clustern</p>	<p>Ein Set von verschachtelten Clustern, welche in einem hierarchischen Baum organisiert werden können.</p>	<ul style="list-style-type: none"> • Distanzmatrix 	<p>Anzahl Dendrogramme mit n Blättern: $\frac{(2n-3)!}{2^{n-2}(n-2)!}$</p> <p>Top-Down (divisiv):</p>	<p>Dendrogramm: Abstand zwischen Clustern: Punkt wo die Cluster verb. werden \rightarrow z.B. 22 !!: Keine Interpretation der Abstände zwischen den Clustern</p>	<pre>d <- dist(x) #Distanzmatrix berechnen hc <- hclust(d, method="ward.D2") # "single" / "complete" / "average" plot(hc) #Dendrogramm erstellen</pre>

<ul style="list-style-type: none"> - Rechenintensiv - Empfindlich gegenüber Ausreissern - Grosse Datensätze schwierig interpretierbar 	<ul style="list-style-type: none"> + Anz. Cluster muss nicht vorgegeben werden + Flexibel in Wahl der Dist.- und Verkn.metriken 	<p>Start mit allen Objekten in einem Cluster, unterteilen in 2 Cluster usw.</p> <p>Bottom-Up (agglomerativ, übliche Form): Start mit allen Objekten als eigenes Cluster, kleinstes Element in Distanzmatrix verbinden, Distanzen in der neuen Gruppe berechnen & wiederholen bis fertig</p> <p>→ Verallgemeinerung für Distanzen zwischen den Clustern: Linkage</p> <ul style="list-style-type: none"> - Single-Linkage: Kleinste Distanz, häufig einz. Abspaltungen - Complete-Linkage: Grösste Distanz, Ausreisser können zu gr. Einfluss - Average-Linkage: Mittlere Distanz, Tendenz zu runderen Clustern - Ward: Varianz zwischen den verbundenen Clustern minimieren, Tendenz zu gleich grossen Clustern 	<p>Maverik: Nur 1 Objekt in einem Cluster Lücken ansehen: hc\$height</p>  <p>Anz. Cluster: hc\$width</p> <p>Heatmap: Visualisieren einer Datenmatrix (blau = Kleine Werte, rot = grosse Werte)</p> <ul style="list-style-type: none"> • Zeilenw. Vergl. um mögliche Cluster zu identifiz. • Farbbalken erlauben das Überprüfen zwischen Clustern und externen kategorialen Variablen 	<pre>rect.hclust(hc, k=3, border="red") #Cluster unterteilen, k vorgeben grhc <- cutree(hc, k=3, h=5) #k=Anzahl Cluster, h=Höhe des Schnittes plot(x, col = grhc, las = 1) #Plotten der Cluster table(Cluster=cutree(res.hc, k=3), diabetes\$class)</pre> <p>Heatmap: str(dat) #Nur numerische Daten dat_ numerisch <- dat[,3:8] #Nur die numerischen Werte auswählen x <- t(as.matrix(dat_ numerisch)) #Als Matrix speichern colnames(x) <- dat\$Car #Namen der Spalten heatmap(x) #Heatmap erstellen, weitere Parameter (distfun, hclustfun, scale) heatmap(as.matrix(x), scale="row")</p> <p>Alternative: library(pheatmap) pheatmap(x, scale="row") #Ergänzen mit Meta-Information, Hinzufügen von Parametern: annotation_col = data.frame(Country=dat\$Country)</p>
<p>DBSCAN Dichtebasiert. Clustering</p> <ul style="list-style-type: none"> - Schwierigkeiten mit Erkennung von Clustern mit untersch. Dichten - Probleme mit hochdimensionalen Daten (Dichte kaum definierbar) 	<ul style="list-style-type: none"> • Distanzmatrix • Datenmatrix <p>+ Für komplexere Daten mit Rauschen und Ausreissern</p> <p>+ Anz. Cluster muss nicht vorgegeben werden</p>	<p>Kernpunkt: Anzahl der Datenpunkte in der ε-Umgebung beträgt mindestens MinPts (inklusive Kernpunkt). Randpunkt: Kein Kernpunkt, liegt aber in der ε-Umgebung des Kernpunkt Rauschpunkt: Weder Kern- noch Randpunkt → Cluster: Kernpunkte plus dazugehörige Randpunkte</p>	<p>kNNdistplot(dist_weapons, k = 19) abline(h=0.15)</p> <pre>res.dbscan <- dbscan(dist_weapons, eps=2200, minPts = 20)</pre> <pre>weapons_london\$cluster <- as.factor(res.dbscan\$cluster)</pre> <pre>plot(x=weapons_london\$Longitude, y=weapons_london\$Latitude, col = res.dbscan\$cluster+1, pch=16)</pre>	<pre>library(dbscan)</pre> <pre>db <- dbscan(df, eps = 0.15, MinPts = 5) #eps = Max. Abstand zw. Pkt.</pre> <pre>par(mfrow=c(1,2), mar=c(3,3,2,1), mgp=c(2, 0.5, 0))</pre> <pre>plot(df[, 1:2], col=multishapes\$shape, las=1, main="multishapes aus factoextra") plot(df[, 1:2], col=db\$cluster+1, las=1, main="DBSCAN")</pre> <p>Parameter: minPts: mind. 3, eps: kann mittels k-Nächste-Nachbarn-Graph abgeschätzt werden (teilw. jedoch auch Domain knowledge)</p>
<p>Modellbasiertes Clustering</p> <p>Maximierung der Log-Likelihood mit Expectation-Maximization-Verfahren (EM-Algorithmus)</p> <p>Der Algorithmus berechnet W'keiten der Zugehörigkeit von Datenpunkten zu Dichten & gruppiert sie basierend in Cluster.</p> <p>- Modell aufgrund von Annahmen über Verteilung der Daten könnte falsch spezifiziert sein</p>	<ul style="list-style-type: none"> • Datenmatrix <p>Daten aus Mischverteilung (multivariat)</p> <p>+ Anz. Cluster muss nicht vorgegeben werden</p> <p>+ Identifikation von Ausreissern</p>	<ul style="list-style-type: none"> - Zufällige Startparameter wählen: Mittelwert μ_c, Kovarianz Σ_c, "Clustergrösse" m_c - Optimales Modell und optimale Anz. Cluster wird durch max. BIC bestimmt, BIC gewichtet Daten mit vielen Parametern starker 	 <p>Linien Plot: Modell und Anzahl Cluster wählen, bei welchem die Linie am höchsten ist.</p> 	<pre>library(mclust)</pre> <pre>mc <- Mclust(Nclus, modelNames="EII") # spezifisches Modell</pre> <pre>mc <- Mclust(Nclus) # Test aller Modelle</pre> <pre>mc\$modelName # bestes Modell</pre> <pre>plot(mc, what="BIC") # Plottet BIC aller Modelle (sofern mc > 1 Modell), kann auch Unsicherheit plotten: what="uncertainty" (je grosser das Symbol desto unsicherere Beobachtung)</pre> <pre>map(mc\$z) # z = Clusterzugehörigkeit</pre> <p>Alternativ: library(factoextra) fviz_mclust(mc, "BIC", palette = "jco") plot(res.mc, what="classification") fviz_mclust(mc, "classification", geom = "point", palette = "jco", xlab="x1", ylab="x2") table(Cluster=res.mc\$classification, diabetes\$class)</p>

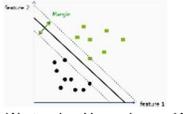
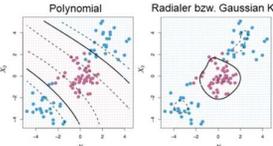
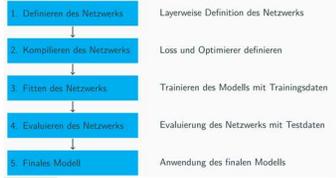
<p>Überwachtes Lernen: Vorbereitung / Allgemein</p>		<p>1. Zufälliges Aufteilen in Trainings- und Test-Daten 2. Anpassen des Modells auf Trainingsdaten (inkl. Preprocessing und Modellselektion) 3. Modellevaluation</p>		<table border="1"> <tr> <td></td> <td>↻ Actual value ↻</td> </tr> <tr> <td rowspan="2">Predicted values</td> <td>TP</td> <td>FP</td> </tr> <tr> <td>FN</td> <td>TN</td> </tr> </table>			↻ Actual value ↻	Predicted values	TP	FP	FN	TN											
	↻ Actual value ↻																						
Predicted values	TP	FP																					
	FN	TN																					
<p>Aufteilung Trainings- / Testdaten (80 % / 20 %) !! Wichtig: Testdaten müssen strikt von Trainingsdaten getrennt werden. Grund: Gefahr, zu optimistisch zu sein. Rückschlüsse der Modellperformance auf neue Daten können nicht gezogen werden. Besser wenn mehr Daten vorhanden, variabler Schätzungsfehler durch zufällige Unterteilung</p> <p>R-Code: library(caret) tindex <- createDataPartition(Studenten\$x, p = 0.8, list = FALSE) # list=TRUE: Liste, list=FALSE: Matrix train_set <- Studenten[tindex,] test_set <- Studenten[-tindex,]</p>	<p>Modellentwicklung <u>Raten:</u> n <- dim(train_set)[1] y_guessing <- sample(c("Male", "Female"), size=n, replace = TRUE)</p> <p><u>Anteil aus Trainingsdaten schätzen:</u> n <- dim(train_set)[1] (p <- sum(train_set\$sex == "Male") / n) y_proportion <- sample(c("Male", "Female"), size=n, prob=c(p, 1-p), replace = TRUE)</p> <p><u>Grenzwert:</u> tapply(train_set\$height, INDEX=train_set\$sex, FUN=function(x) return(c(mean(x), sd(x)))) → untere Grenze: MW - 2*Std. fun_cutoff <- function(x, cutoff) ifelse(x > cutoff, "Male", "Female") y_2sd <- fun_cutoff(train_set\$height, cutoff=157)</p> <p>oder Prüfung verschiedener Werte: cutoff <- seq(150, 200) accuracy <- rep(NA, length(cutoff)) for (i in 1:length(cutoff)){ y_hat <- fun_cutoff(train_set\$height, cutoff=cutoff[i])</p>	<p>Evaluation anhand Testdaten y_acc <- fun_cutoff(test_set\$height, cutoff=164) → Meist tiefer als bei Trainingsdaten mean(y_acc == test_set\$sex)</p> <table border="1"> <thead> <tr> <th>Metrik</th> <th>Calculation</th> <th>Was wird gemessen</th> </tr> </thead> <tbody> <tr> <td>Accuracy</td> <td>$\frac{TP+TN}{n_{test}}$</td> <td>Standard Metrik, unterschiedliche Kosten für Fehler werden nicht berücksichtigt.</td> </tr> <tr> <td>Precision</td> <td>$\frac{TP}{TP+FP}$</td> <td>Wie viele der vorhergesagten ⊕-Objekte sind richtig?</td> </tr> <tr> <td>Recall/Sensitivität</td> <td>$\frac{TP}{TP+FN}$</td> <td>Wie viele der ⊕-Objekte wurden richtig vorhergesagt?</td> </tr> <tr> <td>Spezifität</td> <td>$\frac{TN}{TN+FP}$</td> <td>Wie viele der ⊖-Objekte wurden richtig vorhergesagt?</td> </tr> <tr> <td>F1-Score</td> <td>$\frac{2 \cdot recall \cdot precision}{recall+precision}$</td> <td>Harmonisches Mittel von Recall und Precision</td> </tr> <tr> <td>Kappa</td> <td>$\frac{accuracy - E[accuracy]}{1 - E[accuracy]}$</td> <td>Relative Verbesserung im Vergleich zu einem zufälligen Prädiktor</td> </tr> </tbody> </table>	Metrik	Calculation	Was wird gemessen	Accuracy	$\frac{TP+TN}{n_{test}}$	Standard Metrik, unterschiedliche Kosten für Fehler werden nicht berücksichtigt.	Precision	$\frac{TP}{TP+FP}$	Wie viele der vorhergesagten ⊕-Objekte sind richtig?	Recall/Sensitivität	$\frac{TP}{TP+FN}$	Wie viele der ⊕-Objekte wurden richtig vorhergesagt?	Spezifität	$\frac{TN}{TN+FP}$	Wie viele der ⊖-Objekte wurden richtig vorhergesagt?	F1-Score	$\frac{2 \cdot recall \cdot precision}{recall+precision}$	Harmonisches Mittel von Recall und Precision	Kappa	$\frac{accuracy - E[accuracy]}{1 - E[accuracy]}$	Relative Verbesserung im Vergleich zu einem zufälligen Prädiktor
Metrik	Calculation	Was wird gemessen																					
Accuracy	$\frac{TP+TN}{n_{test}}$	Standard Metrik, unterschiedliche Kosten für Fehler werden nicht berücksichtigt.																					
Precision	$\frac{TP}{TP+FP}$	Wie viele der vorhergesagten ⊕-Objekte sind richtig?																					
Recall/Sensitivität	$\frac{TP}{TP+FN}$	Wie viele der ⊕-Objekte wurden richtig vorhergesagt?																					
Spezifität	$\frac{TN}{TN+FP}$	Wie viele der ⊖-Objekte wurden richtig vorhergesagt?																					
F1-Score	$\frac{2 \cdot recall \cdot precision}{recall+precision}$	Harmonisches Mittel von Recall und Precision																					
Kappa	$\frac{accuracy - E[accuracy]}{1 - E[accuracy]}$	Relative Verbesserung im Vergleich zu einem zufälligen Prädiktor																					
<p>Fehlerarten Fehler Trainingsdaten: Trainingsfehler / In-Sample-Fehler Fehler neue Daten: Verallgem.fehler / Testfehler / Out-of-Sample-Fehler</p> <p>Modellkomplexität mitberücksichtigen (Strafe für Komplexität) → bei 2 Modellen mit ähnlichen Verallgemeinerungsfehlern sollte das einfachere Modell dem komplexeren vorgezogen werden.</p> <p>Overfitting Modell passt zu gut zu den Trainingsdaten, hat aber einen grossen Testfehler / Kleine Genauigkeit.</p> <p>Leave-One-Out Kreuzvalidierung (LOOKV)</p>																							

<ol style="list-style-type: none"> 1. Daten in n Datensätze aufteilen, ein Datensatz weglassen 2. Modell mit allen n Datensätzen trainieren und mit ausgeschlossener Punkt validieren 3. Testfehler aus der Anzahl Fehlklassifizierungen schätzen <p>→ Zufälligkeit nicht mehr gegeben (+) → Kleinster Bias (+) → Numerisch, somit komplex/langsam (-)</p>	<pre>accuracy[1] <- mean(y_hat == train_set\$sex) max(accuracy) cutoff[which.max(accuracy)]</pre> <p>Erste Bewertung: $Accuracy (Richtigkeit) = \frac{Anz. richtig klassifiziert}{Anz. richtig + An. falsch}$ $Errorrate = 1 - Accuracy$ R-Code: mean(y_proportion == train_set\$X)</p>	<p>Konfusionsmatrix (tab <- table(predicted = y_acc, actual = test_set\$sex))</p> <p>mean(y_acc [test_set\$sex == "Male"] == test_set \$X[test_set\$X=="Male"]) # Acc. nach Geschlecht</p> <p>TP + TN / SUM: conf_m[1,1] + conf_m[2,2] / sum(conf_m)</p>
<p>Bootstrap-Methode</p> <ul style="list-style-type: none"> - B neue Stichproben aus den Daten erzeugt - Gleiche Grösse n wie die urspr. Stichproben, Ziehen mit Zurücklegen → Idee: Datensatz als Ersatz für unbekannte zugrundeliegende Verteilung und ziehe daraus neue Stichproben <p>- Bootstrap-Stichprobe aus den Trainingsdaten Validierung auf den Daten, die nicht in der Bootstrap-Stichprobe sind</p>	<p>No-Free-Lunch-Theorem</p> <ul style="list-style-type: none"> - Vereinfachung um überflüssige Details zu eliminieren <p>Gibt kein Modell, dass für jedes Problem am besten funktioniert</p>	<p>K-fache Kreuzvalidierung</p> <ol style="list-style-type: none"> 1. Daten in k sich nicht überschneidende Gruppen einteilen 2. Modell k-mal trainieren (wieder für alle ausser einer Gruppe, jedes Mal diese Gruppe validieren) 3. Testfehler aus der Anz. falsch klassi. Beobachtungen schätzen <p>→ Schneller als LOOKV → Ergebnisse weisen gewisse Zufälligkeit auf → Etwas grösserer Bias als LOOKV, in der Regel jedoch unproblematisch & besser da geringere Varianz als bei LOOKV 10-fache Kreuzvalidierung oft Standard</p>

Überwachtes Lernen: Klassifikation				
Methode	Anwendung / Bemerkungen	+ / -	R-Code / Parameter	
<p>kNN Nächste-Nachbar-Klassifizierer</p>	<ul style="list-style-type: none"> • lokale Strukt: ähnliche Datenpunkte gehören tendenziell zu selben Cluster • nicht linear separierbare Daten Klasse der Mehrheit der nächsten Nachbarn (euklidische Distanz) wird übernommen <p>Bestimmung von k:</p> <ul style="list-style-type: none"> - Kleines k : Zu komplexes Modell - Grosses k : einfaches Modell - Pattsituation verhindern: 2 Klassen: ungerades k wählen > 2 Klassen: k > 1 wählen 	<p>+ Einfach + Flexible Erkennung von Klassen + Gut bei Daten mit lokalen Clustern / Variabilität + Keine Annahme von Verteilung</p> <p>- Empfindlich gegenüber Ausreissern - Rechenintensiv - Wahl k Wert</p> <p>→ nicht linear separierbare Daten</p>	<pre>library(caret) res.knn <- train(y = trainDat\$class, x = trainDat[, -class], method="knn", tuneGrid=data.frame(k=seq(1,20,2)) # Preprocessing: Parameter preprocess = c("center", "scale", "pca") anfügen res.knn <- train(form = class ~., data=trainDat, method="knn", ...) pred <- predict(res.knn, newdata=testDat) #Vorhersage ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 5) method : - "repeatedcv": wiederholte k-fache Kreuzvalidierung - "cv": normale k-fache Kreuzvalidierung - "LOOCV": leave-one-out Kreuzvalidierung - "boot": !!!Bootstrap (Standard)!!! - "none": nichts number: k für k-fach (bei boot: 25, bei KV: 10) repeats: Anz. Wiederholungen der wiederholten Kreuzvalidierung model <- train(x, y, method = "knn", trControl = ctrl,tuneGrid=data.frame(k=1:10)) plot(model) #k bei grösster Accuracy wählen</pre>	<p>Ohne caret: library(class) knn(train=trainDat[, -"label"], cl=trainDat\$label, test=newDat, k=k) confusionMatrix(dat=pred, reference = testDat\$Species)</p> <p>Accuracy bei unbalancierten Daten nicht geeignet</p> <p>Fehlerrate: mean(pred_test != x_train\$Y) Bei k=1: Fehlerrate 0, da nur eigenes Label Klasse definiert</p>
<p>Naive-Bayes-Klassifizierer</p>	<p>Bayes Klassifizierer Beobachtung wird zur wahrscheinlichsten Klasse, basierend auf den Features, zugewiesen. Entscheidungsgrenze: Wo die W'keit für beide Klassen gleichgross ist</p> $P(Y = C_i X) = \frac{Anfangsw'keit Klasse i * W'keit nach Beobachtung von x}{Randw'keit von X, W'keit für Prädiktorvariablen}$ <p>Likelihood bei multivariaten Daten komplex → Naive Bayes-Klassifikator</p> <p>Naiver Bayes Klassifikator Annahmen: Alle Features unabhängig & gleiches Gewicht</p>	<p>+ robust gegen fehlende Daten + Unabh. von Anz. Var.</p> <p>- Annahme, dass alle Var. unabh. voneinander sind - Bei starker korrelation der Variablen, keine guten Ergebnisse</p> <p>→ Gut geeignet für Textdaten, fehlende Daten, grosse Merkm.räume & geringe Stichprobengrösse</p>	<pre>library(naivebayes) ctrl <- trainControl(method = " ", number=x, repeats=y) r.nb <- train(x=einkaufDat[,1:3], y=einkDat\$eink, method="naive_bayes", trControl = trainControl(method="none"), tuneGrid = data.frame(laplace=0, usekernel=FALSE, adjust=0)) # Bei FALSE ist adjust zwingend nötig predict(r.nb, newdata=data.frame(alter="21-30", einkommen="mittel", geschlecht="Frau")) Parameter: Laplace: Faktor für Laplace-Korrektur Usekernel: Empirische Dichteschätzung für numerische Daten Adjust: Bandbreite für Dichteschätzung Alternative:library(e1071) naiveBayes(Einkauf ~., data=dat)</pre>	<p>Problem wenn Häufigkeiten = 0: Laplace Korrektur / Smoothing → Zähler: +konstanter Faktor (meistens zwischen 1 und 2) → Nenner: Faktor mal die Anzahl Features dazuzaddieren</p> <p>Umgang mit stetigen Variablen: Können diskretisiert werden, je nach Implementierung aber auch direkt verwendet werden (sofern normalverteilt, nichtparametr. Dichteschätzung)</p>

<p>LDA / QDA Lineare & quadratische Diskriminanzanalyse</p>	<ol style="list-style-type: none"> Modellierung Verteilung der Features X der verschiedenen Gruppen in den Trainingsdaten individuell. Verwende das Theorem von Bayes um die bedingte Wahrscheinlichkeit $\rightarrow P(Y = C_i X)$ zu bestimmen Ordnung der Beobachtung der Klasse mit der höchsten A-posteriori Wahrscheinlichkeit zu <p>Normalverteilung, Unabhängigkeit vorausgesetzt</p> <ul style="list-style-type: none"> Transformationen können helfen um Daten «normalverteilter» erscheinen zu lassen (log & Wurzel für Exp.verteilungen / Box-Cox für schiefe Verteilungen) Ausreisser entfernen für unverzerrtere Schätzungen Gleiche Varianz: LDA geht davon aus, dass Eigenvariable die gleiche Varianz hat. \rightarrow Standardisieren, sodass MW von 0 und eine Std.abw. von 1 haben 	<p>+ LDA: red. Dimensionalität + Effizient bei kl. Datensatz. + QDA: Berücksichtigung Untersch. Kovarianzstruktur</p> <p>LDA: - Annahme norm. Verteilung – Empfindlich gegenüber Ausreissern - nur für separierbare lineare Daten</p> <p>QDA: - Achtung: Overfitting, Anzahl Var > Anzahl Datpkt</p>	<pre>library(caret) ctrl <- trainControl(method = " ", number=x, repeats=y) # Modellfit LDA lda_fit <- train(y= trainDat\$Species, x = trainDat[, 1:4], method = "lda") res_lda = predict(lda_fit, testDat) # Vorhersage confusionMatrix(data=res_lda, reference = testDat\$Species) # Modellfit QDA qda_fit <- train(y= trainDat\$Species, x = trainDat[, 1:4], method = "qda", ...) res_qda = predict(qda_fit, testDat) confusionMatrix(data=res_qda, reference = testDat\$Species) # Vorhersage</pre>	<p>\rightarrow Gerade gestrichelte Linie deutet auf lineare Diskr.analyse hin, bei quadratischer \rightarrow Gebogene Linie</p> 															
<p>Logistische Regression</p> <p>Multinomiale logistische Regression</p>	<p>W'keit: $\frac{5}{5+3}$ zwischen 0 und 1 Gegenw'keit: $1 - W'keit$ '' Odds: $\frac{5}{3}$ zwischen 0 und ∞ Odds aus W'keiten: $W'keit / Gegenw'keit \rightarrow p / (1 - p)$ log(Odds): zwischen $-\infty$ und ∞</p> <p>$\log\left(\frac{p(Y X)}{1-p(Y X)}\right) = \beta_0 + \beta_1 * x \rightarrow$ Schätzung β_0/β_1 mittels Maximierung Likelihood</p> <p>Multinomiale logistische Regression $\rightarrow > 2$ Klassen Eine der k Klassen wird zur Referenzkategorie k_0 (optimal grösste Klasse) Für alle anderen Klassen wird eine logistische Regression gegenüber der Referenzkat. durchgeführt:</p> <p>$\log\left(\frac{p(Y_i = k X_i)}{p(Y_i = 0 X_i)}\right) = \text{logit}(p_k)$, Total $p_k = 1$</p>	<p>Log Reg: + W'keitsvorhersagen für Klassenzugehörigkeit + Gut für grosse Dat.sätze - Anz. Beobachtungen > Anz. Variablen - bei nicht-linearen Daten schlechte Vorhersagen - Ann. unabh. Beob.</p> <p>Multi Reg: + Mehrklassen-Klassifikation + Modellierung komplexe Zusammenh. (Kl. & Var.) - Ann. lineare Beziehung - Ann. unabh. Beob. - Ausreichend Daten, um robuste Schätzung der Koeffizienten zu erhalten</p>	<p>Logistische Regression</p> <pre>ctrl <- trainControl(method = " ", number=x, repeats=y) library(caret) Bankrott\$Aktiv <- as.factor(Bankrott\$Aktiv) train_index <- createDataPartition(Bankrott\$Aktiv, p = 0.75, list = FALSE) res.glm <- train(y = Bankrott[train_index, 'Aktiv'], x = Bankrott[train_index, -which(names(Bankrott) == 'Aktiv')], method = "glm", family = "binomial") #method = "polr" für ord. Var. summary(res.glm)\$coefficients confusionMatrix(predict(res.glm, newdata=Bankrott[-train_index,]), Bankrott[-train_index,]\$table) Für W'keiten: predict(res_model, newdata, type = 'prob')</pre> <table border="1"> <thead> <tr> <th></th> <th>Estimate</th> <th>Std. Error</th> <th>z value</th> <th>Pr(> z)</th> </tr> </thead> <tbody> <tr> <td>(Intercept)</td> <td>2.3922733</td> <td>1.0093421</td> <td>2.370131</td> <td>0.017781770</td> </tr> <tr> <td>distance</td> <td>-0.6500704</td> <td>0.2296738</td> <td>-2.830407</td> <td>0.004648886</td> </tr> </tbody> </table> <p>Koeffizienten sind immer $\log\text{-Odds} \rightarrow \exp(x\\$coefficients[2])$</p>		Estimate	Std. Error	z value	Pr(> z)	(Intercept)	2.3922733	1.0093421	2.370131	0.017781770	distance	-0.6500704	0.2296738	-2.830407	0.004648886	<p>Multinom. Regression</p> <pre>ctrl <- trainControl(method = "cv", number = 10) multinom_fit <- train(y = iris[train_index, 'Species'], x = iris[train_index, -which(names(iris) == 'Species')], method = "multinom", trControl = ctrl, tuneLength=7, trace = FALSE) confusionMatrix(predict(multinom_fit, newdata=iris[-train_index,]), iris[-train_index, 'Species'])</pre>
	Estimate	Std. Error	z value	Pr(> z)															
(Intercept)	2.3922733	1.0093421	2.370131	0.017781770															
distance	-0.6500704	0.2296738	-2.830407	0.004648886															
<p>Klassifikationsbäume</p>	<p>Intuitives Klassifikationsverfahren um Objekte mittels (binären) «ja/nein»-Fragen einer Klasse zuzuordnen. Abhängig von der zuvor gestellten Frage. Werden in Baumstruktur zusammengetragen.</p> <p>Baum verfügt über folgende Elemente: Wurzel, Knoten für Aufteilung, Äste, Endknoten/Blätter</p>  <p>Trainingsfehler liefert keine gute Schätzung, wie gut der Baum auf ungesehene Daten performen wird Geeignete Komplexität: Fehler = Bias + Varianz + unreduzierbarer Fehler</p> <p>Stopregeln: (Problem, wenn Wachstum zu früh gestoppt wird)</p> <ul style="list-style-type: none"> $\Delta I(N)$ soll Schwellenwert nicht unterschreiten Min. Anz. Objekte, bei welcher ein Knoten aufgeteilt wird Mind. Objekte, welche im Endknoten noch vorhanden sein müssen <p>Besserer Ansatz: Baum vollständig wachsen lassen & zurückschneiden Cost-complexity Pruning: $R(\alpha) = \text{Fehlerrate in den Endknoten} + \alpha * T$ $R(\alpha)$: Mass, welches Modellgüte & Komplexität verbindet T: Anzahl der Endknoten A: Komplexitätsparameter \rightarrow Ziel $R(\alpha)$ so klein wie möglich, finden durch Kreuzvalidierung</p>	<p>+ robust gegenüber Ausreisser + keine Verteilungsann. + keine Variablentransf. notwendig. + Umgang mit kat. & num. Werte + erwartungstr. Schätzer</p> <p>- führt schnell zu Overfitting - Instabilität der Bäume, Interpretation leicht änderbar - benötigt grosse Datenmengen</p>	<pre>ctrl <- trainControl(method = " ", number=x, repeats=y) library(caret), modelLookup("rpart") ctrl <- trainControl(method = "none") r.tree <- train(y = iris[train_index, 'Species'], x = iris[train_index, -which(names(iris) == 'Species')], method = "rpart", trControl = ctrl, tuneGrid = data.frame(cp = 0.01)) # Parameter für minimale Anz. Objekte für Trennung: control = rpart.control(minsplit=20), Mindestanzahl Objekte im Endknoten: control = rpart.control(minbucket = round(minsplit/3)) plot(r.tree\$finalModel, margin=0.2, uniform=TRUE) # Visualisierung text(r.tree\$finalModel) Alternative: library(rpart) tree <- rpart(Species ~., data = my.train)</pre> <p>Kreuzvalidierung cost-complexity Pruning: library(caret) ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3) tree <- train(y = y, x = X, method = "rpart", trControl = ctrl, tuneGrid = data.frame(cp = ...)) plot(tree)</p>	<p>Bestimmung der Verunreinigung (impurity measures): Fehlerrate: $1 - \max(p_i)$, $i = i\text{-te Klasse}$ Gini-Index: Erwartet Fehlerrate beim Knoten N, falls Label zufällig von den am Knoten vorhandenen Klassen selektiert werden: $1 - \sum_{i=1}^g p_i^2$ Entropie: $-\sum_{i=1}^g p_i * \log_2(p_i)$ $\Delta I(N) = I(N) - \alpha_i \cdot I(N_i) - (1 - \alpha_i) \cdot I(N_n) \rightarrow \alpha_i = \text{Ant. Beob. Unterg. } N_i$ \rightarrow Aufteilung so, dass $\Delta I(N)$ maximal</p> <p>Bias-Varianz Tradeoff Bias: Differenz zwischen der mittl. Vorhersage unseres Modells und dem richtigen Wert. \rightarrow hoher Bias: Zu stark vereinfachtes Modell mit hoher Fehlerrate bei Trainings- und Testdaten</p> <p>Varianz: Variabilität der Modellvorhersage für einen bestimmten Datenpunkt \rightarrow Hohe Varianz: Schneiden bei Trainingsdaten sehr gut ab, haben aber hohe Fehlerquoten bei ungesehenen Daten, Andeutung auf zu komplexes Modell</p> <p>Gleichgewicht: Gesamtfehler = Bias² + Varianz + Irreduzierbarer Fehler</p> <ul style="list-style-type: none"> \rightarrow Testen mehrerer Modelle (suche Gleichgew.) \rightarrow Kreuzval. / Anpassung Hyperparameter \rightarrow Regularisierung z.B. mittels Strafterm (Vermeidung Overfitting) \rightarrow Feature engineering (Bearb. Eingabefeatures) \rightarrow Verwendung Ensemble-Methoden 															

Überwachtes Lernen: Ensemble Methoden		Viele Klassifikatoren zusammen verwenden, um leistungsfähiger zu sein → Gesamtklassifizierer (Erfolgreich wenn einz. Klassifizierer im Schnitt unverzerrt) Voraussetzung Einzelkl.: stochastisch unabhängig, Variabilität reduzieren	Stacking: Kombination heterogener schwachen Klassifikatoren unter Verwendung von Metamodellen → Unterteilung in mehrere Trainsätze	
Überwachtes Lernen: Ensemble Methoden		Ensembles sind nur dann besser, wenn jeder Klassifizierer besser ist als zufälliges Raten!	Bagging: Kombination von homogenen, schwachen Klassifikatoren, die parallel anhand unabhängiger Bootstrap-Stichproben trainiert werden (z.B. Random Forest)	
Überwachtes Lernen: Ensemble Methoden		Vorteile: Verbesserung der Genauigkeit, Machen Modell robuster Nachteile: Reduziert Interpretierbarkeit eines Modells, zeitaufwändig	Boosting: Iteratives Verfahren (adaptive Modifikation der Trainingsdaten durch stärkere Fokussierung auf zuvor falsch klassifizierte Datensätze, z.B. Gradient boosting, ADABOOST, XGBoost)	
Methode	Anwendung / Vorgehen	+ / -	R-Code / Parameter	
Random Forest	<ol style="list-style-type: none"> Trainingsset erstellen: N mal mit Zurücklegen aus allen N verfügbaren Trainingsfällen ziehen (Bootstrapsamples) → restliche Daten für Validierung (vorher abtrennen) Ungestutzten Klassifikationsbaum mit Modifikation erstellen: Bei jedem Knoten wird die beste Aufteilung nur unter einer zufälligen Auswahl von mtry Features gesucht, Entkorrelation mit allen Mitteln → Maximierung Varianzreduktion indem Korrelation zwischen Bäumen klein gehalten wird → Verhindert Overfitting Vorhersage aus Mehrheitsabstimmung aller Bäume <p>Kategoriale Variablen mit Random Forest Möglich, kann jedoch komplex werden bei vielen Faktorstufen One-hot-encoding bei Random Forest nicht anwenden!</p>	<ul style="list-style-type: none"> + Robust gegenüber Ausr. + Benchmark Methode + Out-of-Bag Fehlerrate + Umgang mit tausend. Eingangsvar. + neigt nicht zu Overfitting + kann mit unbalancierten Daten umgehen + Auf Reg.probleme anwendbar 	<pre>library(caret) fitC <- trainControl(method = "oob") r.rf <- train(Species ~., data=iris, method = "rf", trControl = fitC, tuneLength=3, ntree=1000) #ntree Anzahl Bäume r.rf\$results r.rf\$finalModel\$confusion #Überprüfung Out of Bag Fehler Mit Ranger: library(caret), library(ranger) rangerOut <- ranger(formula = survived ~ ., data=trainDat, num.trees = 1000, mtry = 2, respect.unordered.factors = "order") #respect.unordered.factors = "order" = target encoding survPred <- predict(rangerOut, data = testDat)\$predictions confusionMatrix(survPred, testDat\$survived) Vergleich verschiedene Modelle: results <- resamples(list(LDA=model_lda, QDA=model_qda, knn=model_knn, RandomForest = model_rf)) summary(results) Visualisierung: bwplot(results) #boxplots dotplot(results) #dot plots</pre>	<p>One-hot-encoding: Erhöht Dimensionalität des Merkmalsraums (kategoriale Variable wird in mehrere Variablen aufgeteilt), kann z.B. bei k-NN nützlich sein</p> <p>Target encoding: Jede Stufe eines kategorialen Merkmals wird durch durchschnittliche Antwort in Zielvariable ersetzt, nützlich in Datenreichem Kontext, Achtung: Mittelwerte können verzerrt werden</p> <p>Umgang mit unbalancierten Daten Alternativer Cutoff wählen & wenn W'keit grösser als dieser Cutoff als Churner klassifizieren.</p> <ul style="list-style-type: none"> - Betrachtung ROC-Kurve (siehe unten) - Anpassung Prior-Wahrscheinlichkeit (siehe unten) <p>Auch möglich durch Sampling Methoden</p>
AdaBoost	<p>Binärer Klassifikator</p> <ul style="list-style-type: none"> - Kombination von schw. Lerner (Bäume häufig nur 1 Stufe) - Nicht alle Modelle gleich gut. Bessere Modelle stärker gewichtet - Jedes Modell berücksichtigt Fehler vom vorgängigen stärker <p>Totaler Fehler des Modell 1 $err_1 = \sum_{i=1}^M w_i \cdot I(y_i \neq M_1(x_i)) = \frac{1}{8}$ Probleme wenn Fehler 0 oder 1: Addieren eines kleinen Fehlerterms</p> <p>Gewicht im Ensemble von Modell 1 $\alpha_1 = \frac{1}{2} \log\left(\frac{1-err_1}{err_1}\right) = 0.97$</p> <p>Beispiele: $w_{23} = \frac{1}{8} \cdot \exp(0.97) = 0.33$ $w_{21} = \frac{1}{8} \cdot \exp(-0.97) = 0.05$ Gewichtung für nächsten Durchlauf, noch unskaliert (muss noch normalisiert werden)</p>		<pre>library(caret), library(adabag) r.Ada <- train(Churn ~., data = churn_train, method = "AdaBoost.M1") confusionMatrix(predict(r.Ada, newdata = churn_test), reference = churn_train\$Churn) → Funktioniert nur bei 2 Klassen, Hyperparameter-Tuning kann sehr lange dauern</pre>	<p>Boosting Allgemein: Anderst als bei Bagging: Die verschiedenen homogenen Klassifizierer werden nicht mehr unabhängig voneinander sondern sequenziell trainiert.</p> <p>Grösstes Gewicht auf Beobachtungen, welche im vorherigen Modell falsch vorhergesagt wurden. → Unterschiedliche Gewichtung der Modelle im Ensemble Hauptaugenmerk auf Reduzierung des Bias → weniger rechenaufwendig</p>
Stochastic Gradient Boosting GBM	<p>Schätzen der Zielvariable mit schwachem Algorithmus (Regression, Klassifikation) Vorhersage schrittweise anhand der Residuen verbessern</p> <p>Vorgehen Gradient Boosting:</p> <ol style="list-style-type: none"> 1. Initialisierung Vorhersagemodells: $\hat{y} = F_0(x)$ 2. Von m = 0 bis M <ol style="list-style-type: none"> a) Berechnung Residuen: $r = y - F_m(x)$ b) Modellieren der Residuen: $\hat{r} = f_m(x)$ c) Neues Modell: $F_{m+1}(x) = F_m(x) + \lambda f_m(x)$ <p>→ Stabilere Optimierung und geringere Gefahr für lokale Optima</p> <p> Hinweis: Stochastic Gradient Boosting ähnelt dem Gradient Boosting, jedoch wird hier bei jedem Schritt nur eine zufällige Teilmenge der Trainingsdaten verwendet. Dies führt zu einer erhöhten Vielfalt und kann das Modell robuster gegenüber Rauschen und Ausreißern machen.</p> <p>M: Anzahl Bäume (1000 – 10000) λ: Lernrate (kleine Rate: reduziert Effekt der einzelnen Bäume, verbessert aber Vorhersage auf lange Sicht. Kleinere Werte liefern bessere Performance – sind jedoch rechenintensiver, Empfehlung 0.01 – 0.001) Allgemeine Gefahr von Overfitting bei Gradient Boosting</p>	<p>Vorgehen Stochastic Gradient Boosting:</p> <ol style="list-style-type: none"> 1. Initialisierung Vorhersagemodells: $\hat{y} = F_0(x)$ 2. Von m = 0 bis M <ol style="list-style-type: none"> a) Berechnung Residuen: $r = y - F_m(x)$ zufälliges Ziehen der Residuen z.B. 50% b) Modellieren der Residuen: $\hat{r} = f_m(x)$ c) Neues Modell: $F_{m+1}(x) = F_m(x) + \lambda f_m(x)$ <p>→ Stabilere Optimierung und geringere Gefahr für lokale Optima</p>	<pre>library(caret), library(gbm), modelLookup("gbm") ctrl <- trainControl(classProbs = TRUE, method = "cv", number = 5, summaryFunction = mnL) #mnLogLoss gbm_fit <- train(Species ~., data = my.train, method = "gbm", trControl = ctrl, distribution = "multinomial", metric = "logLoss", verbose = FALSE, tuneGrid = data.frame(expand.grid(n.trees = c(200, 500, 1000, 2000), interaction.depth = seq(2, 4), shrinkage = 0.01, n.minobsinnode = 10)) # Alternative Verteilungen z.B. "bernoulli" oder "gaussian" confusionMatrix(predict(gbm_fit, newdata = my.test), my.test\$Species)</pre>	<p>Verlustfunktion definiert, was als $\hat{y} = F_0(x) + \sum_{m=1}^M \lambda \cdot f_m(x) = F(x)$ Fehler zwischen den vorhergesagten Werten \hat{y} und den tatsächlichen Werten y angesehen wird. → Quantifiziert, wie gut das Modell $F(x)$ sich der Zielvariable nähert.</p> <ul style="list-style-type: none"> - Vorhersage falsch: Verlustfunktion hat hohen Wert - Variabel anwendbar: z.B. Regression oder Klassifikation - Ziel: Minimierung Verlustfunktion <p>Regression: MSE oder MAE (Mean absolute error)</p> <p>Klassifikation: 0-1 Verlust oder Negative Log Likelihood (Cross-Entropie für W'keitsvorausagen)</p> <p>Regression MSE Klassifikation Log Loss (ohne Herleitung)</p> $L(\hat{y}_i) = (y_i - \hat{y}_i)^2$ $L(\hat{y}_i) = y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)$ $\frac{dL}{d\hat{y}_i} = -2(y_i - \hat{y}_i)$ $\frac{dL}{d\log(\text{odds}_i)} = -y_i + \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}}$ $\frac{dL}{d\hat{y}_i} \propto -r_i$ $\frac{dL}{d\log(\text{odds}_i)} = -y_i - \beta_i = -r_i$

<p>XGBoost</p>	<ul style="list-style-type: none"> - Bietet mehr Regularisierungsoptionen, um Modellkomplexität zu reduzieren und verringert so Overfitting - Berücksichtigt bei Baumentwicklung nur eine reduz. Anzahl Features - Betrachtet Gradient 2. Ordnung -> findet bessere Richtung zur Reduktion Lossfunktion 	<p>Vorteile</p> <ul style="list-style-type: none"> • sehr leistungsfähig • flexibel <p>Nachteil/Gefahren</p> <ul style="list-style-type: none"> • viele Hyperparameter • anfällig auf Overfitting • Interpretierbarkeit <p>XGBoost ist eine Optimierung des Gradient Boosting-Algorithmus</p>	<pre>library(caret), library(xgboost), modelLookup("xgbTree") ctrl <- trainControl(classProbs = TRUE, method = "none", summaryFunction = mnLogLoss) xgb_fit <- train(Species ~ ., data = my.train, method = "xgbTree", trControl = ctrl, metric="logLoss", verbose=FALSE, tuneGrid = data.frame(nrounds = 50, max_depth = 3, eta = 0.3, gamma = 0, colsample_bytree = 0.6, min_child_weight = 1, subsample = 0.5)) # one-hot encoding, nrounds: Anz. Bäume, max_depth: Baumtiefe, eta: Lernrate, gamma: Regularisierung Baumkomplexität, colsample_bytree: Wie viel % der Daten pro Baum, min_child_weight: Mind. Beob., subsample: % der zufälligen Beobachtungen confusionMatrix(predict(xgb_fit, newdata=my.test), my.test\$Species)</pre>	
<p>Support Vector Machine</p>	<p>Flexibles Klassifikationsverfahren mit oft sehr guter Performance – gut geeignet bei vielen Features oder limitierter Anzahl Beobachtungen</p> <ul style="list-style-type: none"> - Jede Beobachtung wird als p-dimensionaler Vektor dargestellt - SVM konstruiert Hyperebene, um 2 Klassen zu trennen  <p>Maximale Margin Klassifikator wählt Ebene so, dass der Abstand von der Hyperebene zum nächsten Trainingspunkt maximal ist.</p> <p>Support Vectors: Trainingsdaten die am nächsten bei der Hyperebene liegen</p> <p>(Abstand zu Hyperebene: $\leq M$), alle anderen Daten tragen nicht zur Festlegung der Grenze bei</p> <p>Optimierungsproblem: $y_i \cdot (\beta_0 + \beta_1 \cdot x_{i1} + \dots + \beta_p \cdot x_{ip}) \geq M$</p>	<p>Oftmals nicht eindeutig trennbare Daten: Fehlklassifizierung erlauben:</p> $y_i \cdot (\beta_0 + \beta_1 \cdot x_{i1} + \dots + \beta_p \cdot x_{ip}) \geq M \cdot (1 - \epsilon_i)$ $\epsilon_i \leq 1/C \rightarrow$ Tuningparameter C darf nicht negativ sein. Grössere Margin erlaubt mehr Fehlklassifikationen, da geringere Kosten. Bei grossen Kosten, kleinere Margin <p>Bei kleinen Margin: Grenzen anfällig falls leichte Veränderungen in Trainingsdaten, Klassifikator mit grosser Variabilität und kleinem Bias, Gefahr von Overfitting</p>	<pre>library(e1071), modelLookup("svmLinear2") Daten als factor definieren ctrl <- trainControl(method="repeatedcv", number=10, repeats=3) model_sva <- train(y~., data=dat, method="svmLinear2", trControl=ctrl, tuneGrid=data.frame(cost=c(0.01, 0.1, 1))) model_sva\$bestTune</pre> <pre>pred <- predict(model_sva, testdat) confusionMatrix(data=pred, reference=testdat\$y)</pre> <p>SVM: Mehr als 2 Klassen → Binärer Klassifikator: Kann nur 2 Klassen trennen</p> <p>Mit K Klassen: OVA (One versus All): Fit K verschiedene SVM Modelle: Jeweils Klasse k gegen den Rest. Neue Beobachtung mit der grössten Distanz zur Hyperebene klassifizieren. OVO (One versus One): Fit alle paarweise SVM-Modelle. Neue Beobachtung in die Klasse klassifizieren, welche die meisten paarweisen Vergleiche gewinnt.</p>	<p>Vorteile</p> <ul style="list-style-type: none"> • Effektive Trennung von Datenpunkte in Klassen • Robust gegenüber Overfitting • Mit OVA oder OVO Ansätze könnte geclustert werden. <p>Nachteile:</p> <ul style="list-style-type: none"> • Sensibel gegenüber Skalierung • Schwierigkeit bei grossen Datensätzen
<p>Support Vector Machine nicht linear (Gaussian Kernel)</p>	<p>Kernel Trick → SVM Polynome manuell hinzufügen, schnell schwierig – daher ist die bessere Lösung Kernels. (Funktion, welche Ähnlichkeit zwischen 2 Beobachtungen quantifiziert.)</p>	<p>Vorteile:</p> <ul style="list-style-type: none"> • Effektiv bei Trennen von nichtlinearen Entscheidungsgrenzen • Robust gegenüber Overfitting • Effektive Verarbeitung von hochdim. Daten <p>Nachteile:</p> <ul style="list-style-type: none"> • Empfindlichkeit gegenüber Skalierung und Hyperparameterwahl • Schwierigkeit bei grossen Datensätzen 	<pre>ctrl <- trainControl(method="repeatedcv", number=10, repeats=3) svm_rad <- train(y~., data=dat, method="svmRadial", trControl=ctrl, tuneGrid=expand.grid(sigma=seq(0.5,5,0.5), C=c(0.01, 0.1, 0.5, 1, 1.5))) svm_rad\$bestTune confusionMatrix(svm_rad, "none")</pre>	<p>Ansatz über Skalarprodukt:</p> <p>Allgemeine Schreibweise mit Kernel $K(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$</p> <ul style="list-style-type: none"> • Linearer Kernel: $K(x_i, x_j) = \sum_{j=1}^p x_{ij} \cdot x_{ij}$ • Polynomieller Kernel d-ten Grades $K(x_i, x_j) = (1 + \sum_{j=1}^p x_{ij} x_{ij})^d$ • Gaussian bzw. Radialer Kernel $K(x_i, x_j) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{ij})^2)$  <p>Beispiele:</p>
<p>Deep-Learning-Ansätze</p> <p>CNN</p> <p>Keras</p>	<p>Deep-Neuronale Netze starten mit Rohdaten und lernen selbstständig während des Trainings, geeignete hierarchische Merkmale zu extrahieren, um diese für die Klassifizierung zu verwenden (End-to-End-Ansatz, bisher Extrahieren von handgefertigten Merkmalen und trainieren eines Modells)</p> <p>Ziel: mathematische Funktionen zu finden die die Attribute auf die Ausgabe abbilden kann</p> <p>Neuronale Netze - Trainieren des Netzes:</p> <ul style="list-style-type: none"> - Gewichte des Netzwerkes sind die Parameter des Modells - Anpassung der Gewichte durch Optimierung der Loss Funktion mit Gradientenverfahren → Iteratives Verfahren - Zufällige Auswahl der Gewichte für die erste Runde - Forward Pass: Input-Werte werden ins Netz gespiesen und durch das Netz propagiert - Backward Pass: Jedes Gewicht und jeder Bias Term wird ein kleines Stückchen in Richtung angepasst, die den Fehler kleiner macht (Gradientverfahren) - Grösste Anpassung über Lernrate (Hyperparameter), hat einen sehr grossen Einfluss <p>Um Overfitting zu vermeiden, ein Teil der Neuronen während dem Training abspalten</p>		<p>Netzwerk Architekturen</p> <p>R-Code: Bilden eines CNNs in keras</p> <pre>model_cnn <- keras_model_sequential() model_cnn %>% layer_conv_2d(filters = 32, kernel_size = c(3, 3), activation = 'relu', input_shape = c(28, 28, 1)) %>% layer_max_pooling_2d(pool_size = c(2, 2)) %>% layer_dropout(rate = 0.25) %>% layer_conv_2d(filters = 64, kernel_size = c(3, 3), activation = 'relu') %>% layer_max_pooling_2d(pool_size = c(2, 2)) %>% layer_dropout(rate = 0.25) %>% layer_flatten() %>% layer_dense(units = 64, activation = 'relu') %>% layer_dropout(rate = 0.5) %>% layer_dense(units = 10, activation = 'softmax') summary(model_cnn)</pre> <p>Deep Learning braucht nicht viel Preprocessing, aber Standardisierung der Daten hilft oft. → !! Aktivierungsfunktion muss zum Wertebereich der Daten passen.</p>	<p>Keras - Workflow:</p>  <p>R-Code: library(keras3), keras3::install_keras()</p> <p># 1. Definieren des Netzwerkes</p> <pre>model <- keras_model_sequential() layer_dense(model, units = 16, input_shape = c(4), activation = 'relu') # input_shape: nur im 1. Layer notwendig, activation: Aktivierungsfunk., rate: Anteil auszulassende Gewichte layer_dropout(model, rate = 0.2) layer_dense(model, units = 16, activation = 'relu') layer_dropout(model, rate = 0.2) layer_dense(model, units = 3, activation = 'softmax') summary(model)</pre>

<p>Netzwerk Architekturen Fully Connected Neuronale Netzwerke (FCNN): für tabellarische Daten, Eingabereihenfolge spielt keine Rolle, geordnete Information geht verloren</p> <p>Convolutional Neural Networks (CNNs): Für Bilder, geteilte Gewichte über das ganze Bild (nicht für jedes Pixel einzelne Gewichte)</p>	<p>Data Augmentation kann helfen um Datenmenge zu erhöhen → leicht modifizierte Kopien von bereits vorhandenen Daten erstellt.</p> <p>Limitierte Daten / Computerressourcen / Zeit:</p> <ul style="list-style-type: none"> - Vortrainierte Modelle, nur letzter Layer wird neu trainiert (Transfer Learning, oder in parallelen Modell Side Learning) - Klassischer Klassifikator (SVM) zur Klassifikation anwenden <p>Bildererkennung: Vortrainierte Modelle anwenden</p>	<p># 2. Kompilieren des Netzwerkes model <- compile(model, optimizer = optimizer_adam(learning_rate = 0.001), loss = 'categorical_crossentropy', metrics = c('accuracy')) <i>Preprocessing der Daten</i></p> <p># 3. Fitten des Netzwerkes history <- fit(model, x=train_scale, y=train_labels, epochs = 100, batch_size = 32, validation_split = 0.2, verbose = 0) plot(history)</p> <p># 4. Evaluieren des Netzwerkes valuate(model, x=test, y=test_labels, verbose=0)</p> <p># 5. Finales Modell predict(model, test[1:10,])</p>
---	--	---

<p>ROC / AUC</p> <p>Güte</p> <p>LogLoss</p>	<p>Sensitivity (True-Positive Rate (tpr)) gegen 1-Specificity (False-Positive-Rate (fpr)) für verschiedene Cut-Off-Werte</p> <p>Modell mit perfekter Trennung geht durch Punkt(1,1) Möglicher Cutoffwert nächster Punkt auf Kurve → ROC-Kurve zeigt, dass Modellvergleiche auf Basis der Accuracy irreführend sein können, da bessere Cutoff-Werte möglich sind.</p> <p>Perfekter Klassifizierer: AUC = 1</p> <p>Normaler Fall: AUC = 0.7</p> <p>Worst Case: AUC = 0.5</p> <p>Modellevaluierung mit AUC Umwandlung von W'keiten zu Klassen → Informationsverlust Modellbewertung anhand der Fläche unter der Kurve (AUC)</p> <p>Effektiv bei unbalancierten Klassen, berücksichtigt Rangfolge, nicht empfindlich gegenüber Absoluten Werten → kann irreführend sein</p>	<p>Veränderung folgender Parameter in predict: pred_prob <- predict(rf_fit, newdata=churn_test, type="prob") pred_alt <- as.factor(pred_prob[,1]>0.2)</p> <p>R-Code: ROC & AUC library(pROC) r_roc1 <- roc(response = Modell1\$obs, predictor = Modell1\$A) # ROC r_roc1\$auc # AUC r_roc2 <- roc(Modell2\$obs, Modell2\$A, , levels=c("B", "A")) # Modell 2 r_roc2\$auc plot(r_roc1, las=1) # Plot → Variable 1 lines(r_roc2\$sensitivities, r_roc2\$specificities, # → Variable 2 col="red", lwd=2, lty=2) legend("bottomright", legend=c("M1", "M2"), lwd=2, lty=c(1,2), col=c("black", "red"), bty="n")</p> <p>Oder direkt in Train: ctrl <- trainControl(classProbs = TRUE, method = "cv", number = 10, summaryFunction = multiClassSummary) iris.knn <- train(y = iris[train_index, 'Species'], x = iris[train_index, -which(names(iris) == 'Species')], method="knn", tControl = ctrl, tuneLength = 3, metric = "AUC")</p>	<p>Modellevaluierung mit Log-Loss Metrik, welche die Güte der Kalibrierung misst (auch Cross-Entropie) $LogLoss = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij})$ N = Anz. Beob., M = Anz. Klassen, y_i: 1 wenn Beobachtung zu Kl. j gehört sonst 0 → Gibt an wie nahe die Vorhersagewahrscheinlichkeit an der tatsächlichen Klasse liegt. Je grösser die Abweichung, desto höher der Log-Loss-Wert → Datenspezifisch → Sinnvoll für W'keitschätzungen, nützlich für Optimierung, Anfällig für Ausreisser</p> <p>R-Code: library(caret) -sum(log(Mod1\$A)*(Mod1\$obs=="A")+log(M1\$B)*(M1\$obs=="B"))/7 caret::mnLogLoss(Modell1, lev = levels(Modell1\$obs))</p> <p>oder direkt in Train: Siehe oben AUC: Parameter metric = "logLoss"</p> <p>Güte einer Wahrscheinlichkeit Kalibrierungskurve: Grafische Beurteilung, ob eine Modell gut kalibriert ist → Wenn vorhergesagte W'keit mit der tatsächlichen Häufigkeit der Ereignisse übereinstimmt, es wird mit den MW gearbeitet</p> <p>R-Code: library(predtools) calibration_plot(data, obs = 'y', pred = 'prob') # Standard 10 Gruppen → Kalibrierung ist perfekt, wenn alle Punkte auf der Winkelhalbierenden liegen, oberhalb = Modell überschätzt W'keiten, unterhalb = Modell unterschätzt W'keiten</p>	<p>ICE-Plots, library(pdp) r_par <- partial(q_fit, pred.var = "Variable", type="classification", which.class="TRUE", prob=TRUE, ice = TRUE) # prob: bei FALSE: logit-Werte, ice: ICE-Kurven, ansonsten nur PDP Kurve plotPartial(r_par, main="ICE Plot für Churn")</p> <p>ICE-Plot kann auch zur Best. der Wichtigkeit angewandt werden: Flache Kurve = Variable nicht wichtig, library(vip) res_vi <- vi(q_fit, train=churn_train, method = "firm", target = "Churn", metric = "accuracy", pred_wrapper = predict , ice = TRUE) vip(res_vi)</p>
--	--	---	--	---

<p>VIP / ICE</p>	<p>Einfluss einer Variable auf Klassifizierung Ind. Conditional Expectation (ICE) Plot / Partial Dependence Plot (PDP)</p> <p>Für jede Beobachtung werden Vorhersagen entlang des Wertebereichs der Variable erstellt. x-Achse: Wertebereich y-Achse: Veränderung W'keit für die Zielklasse (Churn) rote Linie: Mittelwert</p> <p>→ Vorsicht Wechselwirkungen werden nicht berücksichtigt</p>	<p>Wichtigkeit einer Variable</p> <ul style="list-style-type: none"> - Modellspezifische Ansätze: Direkt mit dem Modell verknüpft Nicht für alle Algorithmen verfügbar Nicht über verschiedene Algorithmen vergleichbar - Modellunabhängige Ansätze: Flexibel & für alle Klassifikationsalgorithmen verfügbar z.B. Permutationsansatz von RF (eine Variable weglassen, je grösser die Abnahme der Vorhersagegenauigkeit, desto grösser der Einfluss der Variable) <p>→ Wenn Modell schlecht, Variable Importance nicht aussagekr.</p>	<p>Modellspezifischer Ansatz mit RF iris_rf <- train(Species ~., data=iris, method="rf", ..., importance = TRUE) varImp(iris_rf, scale = FALSE) # für Wichtigkeit einer Variable</p> <p>Modellunabhängiger Ansatz: library(vip) r_vi <- vi(q, train = churn_test, target = "Churn", #q = train(...) method = "permute", metric="accuracy", pred_wrapper = predict, nsim = 3) # nsim <i>Anzahl Wiederholungen (Zufallsprozess)</i> vip(r_vi, geom = "boxplot")</p>	<p>Umgang mit fehlenden Werten</p> <ul style="list-style-type: none"> - Missing completely at Random (MCAR): Kein Zusammenhang zwischen den fehlenden Daten und den restlichen Werten (zufällige Teilmenge) - Missing at Random (MAR): Fehlende Werte haben einen Zusammenhang mit den beobachteten Werten, aber nicht mit den fehlenden Werten selbst - Missing not at Random (MNAR): Beziehung zwischen fehlendem Wert und dessen Wert <p>Ansätze: Complete Case Analysis (MCAR, MAR – bei MNAR !! → Bias): - Entfernen von Beobachtungen (Zeilen) train(..., na.action = na.omit) - Wenn z.B. > 60 % eines Features fehlt evtl. die ganze Spalte entfernen</p>
-------------------------	--	--	--	--

<p>Sampling Methoden</p> <p>Prior-W'keit</p> <p>Feature Engin.</p> <p>Fehl. Werte</p>	<p>Sampling Methoden Prior Sampling: Falls Dysbalance bereits zu Beginn bekannt, kann die Datenaufnahme angepasst werden</p> <p>Up-Sampling: Ziehen aus der kleineren Klasse bis Klassen ausgeglichen sind (SMOTE)</p> <p>Down-Sampling: Beobachtungen aus Mehrheitsklasse löschen / auf Anzahl der Minderheitsklasse anpassen (!Verlust von Daten) → Erst nach dem Aufteilen in Trainings- und Testdaten durchführen</p> <p>R-Code: Down-Sampling bei Random Forest ctrl <- trainControl(method = "oob") rf_fit <- train(x = x, y = y, data = dat, method = "rf", ..., strata=dat\$y, sampszie=c(rep(min(table(dat\$y)), nlevels(dat\$y))) data_down <- downSample(x=diabetes[train,], y=diabetes\$diabetes[train], list=FALSE)</p>	<p>Anpassung Prior-Wahrscheinlichkeit Bei einigen Klassifizierern kann eine Prior-Wahrscheinlichkeit eingestellt werden: Naive Bayes-Klassifikator / Diskriminanzanalyse Standardmässig wird diese aus der Verteilung in den Daten bestimmt, kann aber auch angepasst werden: Dafür muss dem train Modell folgender Parameter mitgegeben werden: lda_fit <- train(Churn ~., method = "lda", prior=c(0.5,0.5), ...)</p> <p>PCA & Eigenwerte / Eigenwertvektoren: Rotationsmatrix A kann durch Eigenvektoren und Kovarianz-Matrix bestimmt werden. Zugehörige Eigenwerte entsprechen den Varianzen der einzelnen Hauptkomponenten</p>	<p>Feature Engineering Umwandlung von Rohdaten in Features (Datenverständnis notwendig) → Oftmals grösseren Einfluss als Algorithmus</p> <p>Textdaten</p> <ul style="list-style-type: none"> - Tokenisierung - Stoppwörter (uninformative Wörter) entfernen - Stemming / Lemmatisierung (Wörter in Stammform bringen) - Document-Term Matrix (Tabelle mit Häufigkeiten der Wörter) - Bag-of-Words (Vektor mit Anz. Beob. für jedes Wort in Wortliste) <p>Bilddaten</p> <ul style="list-style-type: none"> - Deep Learning Methoden - Features aus Bildern extrahieren (Features sind Vektoren) 	<p>Umgang mit fehlenden Werten</p> <ul style="list-style-type: none"> - Missing completely at Random (MCAR): Kein Zusammenhang zwischen den fehlenden Daten und den restlichen Werten (zufällige Teilmenge) - Missing at Random (MAR): Fehlende Werte haben einen Zusammenhang mit den beobachteten Werten, aber nicht mit den fehlenden Werten selbst - Missing not at Random (MNAR): Beziehung zwischen fehlendem Wert und dessen Wert <p>Ansätze: Complete Case Analysis (MCAR, MAR – bei MNAR !! → Bias): - Entfernen von Beobachtungen (Zeilen) train(..., na.action = na.omit) - Wenn z.B. > 60 % eines Features fehlt evtl. die ganze Spalte entfernen</p>
---	---	--	--	--

		<p>Einschub: Eigenwerte und Eigenvektoren</p> $M \cdot x = \lambda \cdot x \quad \begin{bmatrix} m_{11} & \dots & m_{1p} \\ \vdots & \ddots & \vdots \\ m_{p1} & \dots & m_{pp} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$ <p>Eigenvektor x einer quadratischen Matrix ist ein vom Nullvektor verschiedener Vektor, dessen Richtung durch Multiplikation mit der Matrix nicht verändert wird. Eigenvektor wird nur gestreckt. Streckungsfaktor λ ist Eigenwert der Matrix.</p>		<p>Imputation von fehlenden Daten (vervollständigen):</p> <ul style="list-style-type: none"> - Einsetzen eines konstanten Wertes der sich von anderen Werten unterscheidet (z.B. neue Faktorstufe) - Einsetzen von Mittelwert, Median oder Moduswert (oftmals nicht empfehlenswert) - Schätzen des Wertes mit prädiktivem Modell - Mehrfache Imputation <p>R-Code: library(caret) knn_imp <- preProcess(iris_miss, method= "knnImpute" / "bagImpute" / "medianImpute") # danach predict(knn_imp, data_missing) Multiple imputation mit library(mice)</p>
--	--	--	--	---

Falsch	In Spezialfällen richtig	Richtig
Um das 90% Quantil von hochdimensionalen Daten zu berechnen, kann man die Beobachtungen der Grösse nach ordnen und dann bestimmt die Beobachtung das 90% Quantil, für die 90% aller anderen Beobachtungen kleiner-gleich sind	Wenn alle Variablen numerisch sind, dann macht es keinen Sinn, für eine PCA die Variablen zu skalieren	Der Jaccard Koeffizient der beiden Punkten A=(0, 1, 0, 0) und B=(0, 1, 0, 0) ist 1
Bei einer 10-fold cross-validation, sind 10% der Daten im Testdatensatz.	Bei einem Lasso-Modell werden die Koeffizienten erwartungstreu geschätzt.	Wenn zwei Beobachtungseinheiten bei einem Klassifikationsbaum in derselben Partition des Feature Raums landen, dann haben sie die gleichen Feature Werte
Die Support Vector Machine ist ein probabilistisches Klassifikationsmodell.		Die PCA besteht aus einer Verschiebung und einer Rotation des Koordinatensystems.
Um einen Clustering-Algorithmus anzuwenden, braucht man eine Zielvariable.		Eine Dimensionsreduktion mit linearer MDS und euklidischen Distanzen führt zum selben Ergebnis wie eine klassische PCA
Bei einer PCA kann die zweite Dimension/Hauptkomponente mehr Varianz erklären als die erste Dimension.		Ist bei t-SNE der perplexity Parameter zu hoch, können Strukturen sichtbar werden, die in den Daten gar nicht vorhanden sind.
Die einzelnen Hauptkomponenten einer PCA können korreliert sein		PAM hat im Gegensatz zu K-Means immer eine Beobachtung als Center für jedes Cluster.
Eine MDS kann mit einer Korrelationsmatrix als Input umgehen		Hierarchische Clustering-Algorithmen benötigen eine Distanzmatrix, um verschachtelte Cluster zu bilden.
Clustering mit K-Means findet immer eine globale Lösung (d.h. globales Minimum).		Bei einer Kreuzvalidierung würden wir, wenn die Rechenleistung keine Rolle spielen würde, immer eine Leave-One-Out Kreuzvalidierung machen.
Wenn ein Punkt nach dem Clustering im Silhouetten-Plot einen negativen Silhouettenkoeffizienten aufweist, heisst das, dass dieser Punkt ein Ausreisser ist und keinem Cluster zugewiesen werden kann.		Ein Vorteil von Ensemble-Verfahren aus schwächeren Lernern ist, dass die Varianz im Ensemble deutlich kleiner wird, als sie es von den einzelnen Lernern war.
Ein knn-Klassifizierer wird dasselbe Ergebnis liefern, egal ob die Daten skaliert oder unskaliert verwendet werden.		
Je komplexer ein Klassifikationsalgorithmus ist, umso kleiner wird sowohl der Trainingsfehler als auch der Testfehler.		