

Unüberwachtes Lernen: Dimensionsreduktion

Zweck Dim.reduktion	Quantifizierung	Distanzmass d $d \leftarrow \text{dist}(\text{dat})$	Euklidische Dist.	Manhattan Distanz	Kategoriell
<ul style="list-style-type: none"> Visualisierung wichtig. Bei hochdim. Daten Visualisierung nicht möglich → Verfahren, um Zusammenhänge/Strukturen zu sehen. Preprocessing für weitere Analysen (bspw. für überwachtes Lernen) Ausreiser detektieren Hilft irrelevante Variablen zu eliminieren oder Noise zu reduzieren 	<ul style="list-style-type: none"> Quantifizierung Ähnlichkeit zwischen zwei Objekten mit Ähnlichkeits-/Distanzmass. Ähnlichkeitsmass s: quantifiziert Übereinstimmung. Je grösser Wert, desto grösser Ähnlichkeit. Distanzmass d: (z.B. euklid. Dist zwischen Punkten). Je kleiner Wert, desto grösser Ähnlichkeit. → Grundlagen für MDS, T-SNE, UMAP und Clustering 	<ul style="list-style-type: none"> O_1, O_2 sind 2 Objekte. Distanz O_1 zu O_2 ist \mathbb{R} Zahl $d(O_1, O_2)$ Distanzmass $d(O_1, O_2)$ → reellwertige Distanzfunktion Eigenschaften <ul style="list-style-type: none"> $d(O_1, O_2) \geq 0$ nicht-negative Funktion $d(O_1, O_1) = 0, d(O_2, O_2) = 0$ Distanz zu sich selbst ist 0 $d(O_1, O_2) = d(O_2, O_1)$ symmetrisch Metrische Distanzen (Metrik) erfüllen zusätzlich: <ul style="list-style-type: none"> $d(O_1, O_2) = 0$ genau dann, wenn $O_1 = O_2$ $d(O_1, O_2) \leq d(O_1, O_3) + d(O_3, O_2)$ Delta-Ungleichung 	<ul style="list-style-type: none"> L2-Metrik → euclidean Bei Beobachtungen im Raum gehen wir intuitiv euklid. Dist. aus. Euklid. Dist. 2 Beobachtungen o_i, o_j, mit p numerisch Features: $d_{ij}(o_i, o_j) = \sqrt{\sum_{k=1}^p (o_{ik} - o_{jk})^2}$ 	<ul style="list-style-type: none"> L1-Metrik → manhattan Ein Block ist eine Einheit. Wie viele Blocks muss man von O_1 nach O_2 zurücklegen. Manhattan Metrik: $d_{ij}(o_i, o_j) = \sum_{k=1}^n o_{ik} - o_{jk}$ Manhattan für hochdimensionale Daten besser, da grosse Abstände nicht quadriert werden. 	<ul style="list-style-type: none"> Objekte haben p nominale Merkmale: Sneath-Übereinstimmungskoeffizient $s = \#Übereinstimmung/p \rightarrow$ factor Ordinale Kategorien mit K Stufen: Transformieren $z_{O_i} = \frac{x_{O_i} - 1}{K - 1} \in [0, 1]$ ordratio
Distanzmatrix D			<ul style="list-style-type: none"> Paarweise Distanz in symmetrischer Matrix. Diagonale sind alle 0. 	Minkowski Distanz (Allgm) $d_{ij}(o_i, o_j) = (\sum_{k=1}^p o_{ik} - o_{jk} ^r)^{\frac{1}{r}}$	Mix Variablentyp
Distanzmasse für binäre Daten					
<ul style="list-style-type: none"> Balancierte Variable: Simple-Matching/Canberra-Metrik symm Unbalancierte Variable: Jaccard Koeffizient/Metrik → asymm 					

`library(cluster); dist <- daisy(dat, metric = "gower", stand = F, type = list(asymm = 3, ordratio = c(7,8))); str(d) #metric = Distanzmass, stand = Standardisierung, type = Behandlung Variable Nr. X als`

Manifold-Hypothese: Reale hochdimensionale Daten liegen auf einer niedrigdimensionalen Mannigfaltigkeiten, die in den hochdimensionalen Raum eingebettet sind.

M	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
Hauptkomponenten Analyse PCA	<ul style="list-style-type: none"> Dimensionsreduktion und Visualisierung hochdim. Daten Methode zur Detektion von Ausreissern Idee: Finde diejenigen Achsen, welche die multidimensionalen Daten am besten erklären. Drehung zentrierten Originalvariablen in Hauptkomponentensystem, so dass diese unkorreliert sind und absteigender Weise möglichst grossen Teil Totalvarianz abdecken. Ziel: Dimensionsreduktion: Verwenden nur die ersten paar Hauptkomponenten und hoffen, dass die Daten dadurch möglichst gut beschrieben werden Faustregel: ~ 80 % der Var_{total} soll durch ersten PC's beschrieben werden (sonst zu wenig aussagekräftig) 	<ul style="list-style-type: none"> Datenmatrix Numerische Variablen Alternative: <ul style="list-style-type: none"> Standardisierung meistens mit <code>mean()</code> und <code>sd()</code> Häufig wäre es besser/sicherer Standardisierung robust durchzuführen. Statt Mittelwert/Standardabweichung robuste Masse Median/MAD nehmen. → Ausreisser in einzelnen Variablen werden besser sichtbar Standardisierung muss für <code>prcomp()</code> manuell durchgeführt werden und PCA anschliessend mit manuell standardisierten Daten durchgeführt werden 	<ul style="list-style-type: none"> + Reduziert Dimensionen + keine Information gehen mit PCA verloren! - Nur für numerische Variablen - Varianzen bzw. Korrelationen sind sensitiv bezüglich Ausreisser → Ergebnis der PCA kann von Ausreissern stark verfälscht werden! - Oftmals genügen die ersten zwei bis drei Hauptkomponenten nicht, um genügend der Var_{total} zu beschreiben (mindestens 80 %). 	<pre>pca <- prcomp(dat, scale = T/F) summary(pca) plot(PC2 ~ PC1, data = pca\$x, xlab = "PC1", ylab = "PC2", col = dat\$Class) var <- pca\$sdev^2 # 80 % Faustregel var_cum <- cumsum(var)/sum(var) plot(1:length(var), var_cum, las = 1) abline(h = 0.8, col = "red", lty = 2) (pc_number <- min(which(var_cum > 0.8))) #Alternative library(ggfortify); biplot(pca) autoplot(pca, colour = 1:3[cancer\$Level])</pre>
Robuste PCA	<ul style="list-style-type: none"> PCA mit robusten Schätzern (robuste Kovarianzmatrix) Ziel Detektion von Ausreissern <ul style="list-style-type: none"> Ausreisser PCA-Raum: Grosse Distanz $SD_i = \sqrt{\sum_{j=1}^k \frac{y_{ij}^2}{\lambda_j}}$ zum Ursprung der PCA → Score-Distanz Orthogonal Ausreisser: nicht richtig erfasst in PCA Rm. Orthogonale Distanz $OD_i = x_i - \hat{\mu} - P \cdot y_i^T$ 	<ul style="list-style-type: none"> Datenmatrix Numerische Variablen 	<ul style="list-style-type: none"> + Reduziert Dimensionen + Nicht ausreiseranfällig - Nur für numerische Variablen 	<pre>library("rrcov") #k = Ausgabe Dimension rob <- PcaHubert(dat, k=2, scale=T) plot(rob@scores[,1], rob@scores[,2], col = c("red", "blue")[as.factor(card\$Class)], pch = 20, xlab = "PC1", ylab = "PC2") plot(x = rob@od); abline(h = rob@cutoff_od)</pre>
Multi-Dimensional Scaling (MDS)	<ul style="list-style-type: none"> Verwendet Distanzen zwischen Beobachtungen und erzeugt Projektion der Punkte in 2D. Wir haben Distanzmatrix von vielen Beobachtungen → Visualisierung in 2-(3)-Dimensionen Klassische metrische MDS: Minimierung Kostenfunktion $Cost = \sum_{i < j} (d_{ij} - d_{ij}^*)^2$ Weitere MDS-Versionen: Ordinale oder nicht-metrische MDS <code>isoMDS()</code>: versucht low-dimensionale Repräsentation der Ränge v. Distanzen zu finden. Cost-Funktion verlangt explizit keine euklidischen Distanzen. 	<ul style="list-style-type: none"> Distanzmatrix Geht für alle Variablentypen (kategoriale → gower) Gibt metrische und nicht-metrische MDS <p>PCA und metrische MDS sind äquivalent, wenn euklidische Distanzen verwendet wurden.</p>	<ul style="list-style-type: none"> + Geht für kategoriale Variablen + Benötigt nur Distanzmatrix + Metrische MDS ist (meistens) auch okay für nicht-euklidische Distanzen. - <code>isoMDS()</code> kann nicht mit Duplikaten (also Distanz 0) umgehen 	<pre># Metrische MDS cmd <- cmdscale(dist, k = 2) plot(cmd, pch = "") #k=Ausgabe Dimension text(cmd, labels = rownames(cmd)) # Nicht-metrische MDS library(MASS); library(cluster) d <- as.matrix(daisy(dat, metric = "gower")) iso <- isoMDS(d, k = 2) #k = Ausgabe Dim. plot(iso\$points, pch = ""); text(iso\$points, labels = rownames(iso\$points))</pre>

Problem mit PCA / metrischem MDS: Gibt fast keinen Grund, dass Daten wirklich in Hyperebene liegen sollten. → Deshalb neue Ansätze (t-SNE/UMAP): Lokale Distanzen korrekt wiedergegeben.

M	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
<p>Stochastic Neighbour Embedding (t-SNE) von Laurens van der Maaten (2008)</p>	<ul style="list-style-type: none"> Versucht Nachbarschaftsbeziehungen im hochdimensionalen Raum zu erhalten und niedriger Dim. darzustellen. Berechnung euklid. Abstand $x_i, x_j: \ x_i - x_j\ ^2 = \sqrt{(x_i^x - x_j^x)^2 + (x_i^y - x_j^y)^2}$ Dichtewert der Distanz $\ x_i - x_j\ ^2 = p_{x_j x_i} = p_{j i}$ unskalierte Ähnlichkeit von x_j zu x_i betrachtet von x_i Skalieren Ähnlichkeitsmasse $p_{j i}$ durch Dividieren der Summe aller Ähnlichkeitsmasse: $p_{j i} = \frac{p(x_j x_i)}{\sum_{k \neq i} p(x_k x_i)}$, somit gilt $\sum_j p_{j i} = 1$ (wobei $p(i i) = 0$) Zwei Punkte haben zueinander unterschiedliche bedingte Wahrscheinlichkeiten $p_{i j}, p_{j i}$ (stammen aus zwei versch. Verteilungen) → $p_{ij} = \frac{p(j i)+p(i j)}{2n}$, $n = \#$ Dimensionen Bedingte Wahrscheinlichkeit: $p_{j i} = \frac{\exp(-\ x_i - x_j\ ^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\ x_i - x_k\ ^2 / 2\sigma_i^2)}$ σ_i ist abhängig vom Punkt x_i und basiert auf Dichte in Nachbarschaft vom Punkt. Varianz wird indirekt über Wahl der Perplexität (Hyperparameter, vereinfacht = Anzahl Nachbarschaftspunkte) definiert: $\text{Perplexity} = 2^{H(P_i)} = 2^{-\sum_j p_{j i} \log_2 p_{j i}}$ (mit steigender Perplexity nimmt σ zu) Entropie $H(P_i) = -\sum_j p_{j i} \log_2 p_{j i} = \sum_j p_{j i} \log_2 \left(\frac{1}{p_{j i}}\right)$ Mass für die erwartete Überraschung (Werte zwischen 0 (keine Überraschung) und $\log_2(\# \text{ Punkte})$ maximale Überraschung). $\log_2 \left(\frac{1}{p_{j i}}\right)$ = Mass für Überraschung, wenn $p_{j i}$ gross → wenig Überraschung σ^2 eher klein → Range für Anzahl lokaler Nachbarn eher klein → kleine Perplexity Wird im niedrig-dimensionalen Raum Gaussverteilung als Ähnlichkeitsmass verwendet, gibt es Crowding-Problem (durch Dimensionsreduktion können nicht mehr alle Strukturen/Distanzen korrekt wiedergegeben werden). t-SNE löst Problem, indem Ähnlichkeiten q_{ij} über t-Verteilung (df = 1) berechnen: $q_{j i} = \frac{(1 + \ y_i - y_j\ ^2)^{-1}}{\sum_{k \neq i} (1 + \ y_i - y_k\ ^2)^{-1}}$ 	<ul style="list-style-type: none"> Daten- oder Distanzmatrix Kategoriale Variablen können nur via <code>gower()</code> Distanzmatrix der Funktion <code>Rtsne()</code> übergeben werden. Datenmatrix geht nur für numerische Variablen. Algorithmus Initialisierung stoch. Testen: <code>set.seed()</code> Berechnung euklid. Abstand zwischen allen Punkten Berechnung Dichtewert Distanz $p_{j i}$ Skalieren Ähnlichkeitsmasse, $p_{ij} = \frac{p(j i)+p(i j)}{2n}$ Berechnung bedingte Wahrscheinlichkeiten Erstellung niedrig-dimensionalen Raum mit gleichen Punkteanzahl wie im ursprünglichen Raum, welche zufällig verteilt werden. Berechnung Wahrscheinlichkeiten q_{ij} (Ähnlichkeiten) im niedrig-dimensionalen Raum (Bsp. 2D). Punkte mit Optimierungsverfahren so anpassen, dass sie möglichst nah an bedingte Wahr'keit p_{ij} im hoch-dimensionalen Raum kommen. Wird durch Minimieren der Kullback-Leibler divergence Verlustfunktion getätigt. Algorithmus mit verschiedenen Seeds durchführen. Zufällige Initialisierung der Optimierung kann bei verschiedenen Durchläufen zu unterschiedlichen Lösungen (lokalen Optima) führen. → Mehrfaches Ausführen des Algorithmus (beste Visualisierung; niedrigster Wert für die Zielfunktion <code>itercosts</code>) Mit Hyperparameter perplexity experimentieren, verschiedene Werte nehmen: <ul style="list-style-type: none"> Je mehr und dichter Daten, desto höher der Wert → Faustregel \sqrt{N} Es gilt: $3 \cdot \text{perplexity} < \text{nrow}(X) - 1$ Hyperparameter Anzahl Iteration (max_int) sollte ausreichend gross gewählt werden für stabile Konfiguration Algorithmus passt Abstände an lokale Unterschiede in Dichte an. So werden dichte Cluster vergrößert und breite Cluster verkleinert. 	<p>+ Lokale Nachbarschaften bleiben eher erhalten</p> <p>- kann Clustering liefern, welche eigentlich gar nicht existiert</p> <p>- Anwendung auf neue Punkte strandrömässig nicht möglich</p> <p>- Clustergrößen bedeuten nichts</p> <p>- Abstände bedeuten möglicherweise nichts</p> <p>- Abstände zwischen Clustern können nicht interpretiert werden</p> <p>- Dimensionsreduktion mit t-SNE als Input für ein Clusterverfahren ist gefährlich, da das Clusterverfahren die Struktur von t-SNE stark übernimmt.</p>	<pre>library(Rtsne) set.seed(7) tsne <- Rtsne(dat, dims = 2, is_distance = F, perplexity = 30, max_iter = 1000) dims: Ausgabedimension [2] pca: Angabe, ob zuerst PCA auf Daten angewendet werden soll [TRUE] pca_center → Datenzentrierung vor PCA [TRUE] pca_scale: Skalierung Daten vor PCA [FALSE] initial_dims: Anzahl Dimensionen, die von PCA verwendet werden [50] is_distance: Indikator, ob X Distanzmatrix ist [FALSE] check_duplicates: Überprüfung, ob Duplikate vorhanden [TRUE] theta: Geschwindigkeit/Genauigkeitskompromiss [0.5] 0.5 schnelle Implementierung; 0 exakte Berechnung perplexity: Hyperparameter; entspricht ungefähre Anzahl erwarteter Nachbarn [30] → Faustregel \sqrt{N}: (5,30), (5,50) max_iter: Hyperparameter; Anzahl der Iterationen [1000] plot_data <- as.data.frame(tsne\$Y) plot_data <- cbind(plot_data, Meta23) ggplot(plot_data, aes(x = V1, y = V2, label = Name, colour = Fraktion)) Dimensionsreduktion als Input für ein Clusterverfahren library(Rtsne); set.seed(60) tsne <- Rtsne(dat, perplexity = 100, max_int = 5000) d_tsne <- as.data.frame(tsne\$Y) names(d_tsne) <- c("tsne1", "tsne2") d_tsne\$label <- as.factor(cancer\$Level) library(dbSCAN) (res.dbSCAN <- dbSCAN(tsne\$Y, eps = 3, minPts = 25)) library(ggplot2) ggplot(data = d_tsne, aes(x = tsne1, y = tsne2, col = res.dbSCAN\$cluster, label = label)) + geom_text(size=3)</pre>

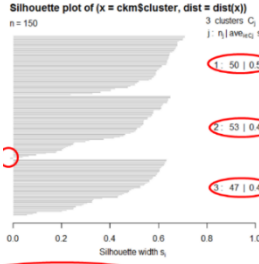
M	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
Uniform Manifold Approximation and Projection (UMAP)	<ul style="list-style-type: none"> Initialisierung fix (spektrales Embedding): Zufälligkeit durch stochastischer Gradientenabstieg (<code>random_state</code>) Hauptunterschiede UMAP zu t-SNE: Abstand im hochdimensionalen Raum wird über Exponentialverteilung bestimmt, (ohne Normalisierung) Mittelung Ähnlichkeiten in beide Richtungen: $p_{ij} = p_{ij} + p_{ji} - p_{ij}p_{ji}$ Der Parameter <code>n_neighbors</code> (Anzahl nächsten Nachbarn) steuert, wie UMAP zwischen lokaler und globaler Struktur abwägt. Bei niedrigen Werten konzentriert sich UMAP mehr auf lokale Struktur, bei hohen Werten auf Struktur Gesamtbild. Ähnlich wie bei <code>perplexity</code> in t-SNE. Initialisierung der Punkte in niedrigen Dimension erfolgt nicht zufällig, sondern über spektrale Einbettung. Für Ähnlichkeit in niedrigen Dimension wird eine Verteilung ähnlich der t-Verteilung verwendet (gesteuert über den Parameter <code>min_dist</code> → bei niedrigen Werten werden die Punkte dichter zusammengepackt). Kostenfunktion mit binärer Kreuzentropie anstelle der Kullback-Leibler-Divergenz Hyperparameter <code>n_neighbors</code> und <code>min_dist</code> wichtig → testen versch. Konfigurationen! 	<ul style="list-style-type: none"> Daten- oder Distanzmatrix Kategoriale Variablen können nur via <code>gower()</code> Distanzmatrix der Funktion <code>umap()</code> übergeben werden. Datenmatrix geht nur für numerische Variablen. Ähnlichkeit hochdimensionalen Raum Berechnung Abstand zwischen Punkten in hoher Dim. Für jeden Punkt wird Abstand unter exponentiell abfallenden Kurve aufgetragen. Kurvenform hängt von Anzahl Nachbarn (einschließlich des Punktes selbst) ab. → $\log_2(\# neighbors)$ entspricht Summe unskalierten Ähnlichkeiten (Dichte). Ähnlichkeiten werden für alle Punktkombinationen auf gleiche Weise erzeugt. UMAP Konzept Nach Initialisierung (Spektral Embedding) werden Punkte so verschoben, dass Ähnlichkeiten in niedrigen Dimension mit Ähnlichkeiten in hohen Dimension übereinstimmen. Die Ähnlichkeit in der niedrigen Dimension wird anhand einer Kurve ähnlich der t-Verteilung bestimmt. Bei der Anpassung wird die Kostenfunktion (binäre Kreuzentropie) minimiert. Hierfür wird stochastische gradient descent Algorithmus verwendet. Dieser hat Zufallskomponente, weshalb mehrere Durchläufe des Algorithmus, trotz eindeutiger Initialisierung, zu unterschiedlichen Ergebnissen führen können. 	<p>Vorteile gegenüber t-SNE</p> <ul style="list-style-type: none"> + Anwendung auf neue Punkte möglich (wenn Datenmatrix) + Dimensionsreduktion auch möglich für mehr als 2-3 Dimensionen + UMAP kann auf Rohdaten ohne PCA-Vorverarbeitung ausgeführt werden. Die PCA kann aber auch hier hilfreich sein. + Anwendung direkt möglich auf grosse Datensätze (performanter, schneller als t-SNE) + Grösserer Fokus auf globale Struktur (Abstände zwischen Gruppen teilweise interpretierbar) als bei t-SNE + Cluster oft besser getrennt <p>– Wie t-SNE kann auch UMAP falsch interpretiert werden.</p> <p>– Clustergrößen in Diagramm haben keine Bedeutung (lokale Abstände sind entscheidend)</p> <p>– Abstände zwischen Clustern haben möglicherweise keine Bedeutung und sind nicht immer sinnvoll</p>	<pre>library(umap) res_umap <- umap(d = iris, n_neighbors = 25, min_dist = 0.05, input = "dist", random_state = 8) plot(res_umap\$layout, col = iris\$Species) n_neighbors: Anzahl nächsten Nachbarn [15] min_dist: Abstand Punkte in Visualisierung [0.1] input: Datenmatrix "data" oder Distanzmatrix "dist" (as.matrix(dist)) random_state: Seed für Realisierung 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) + scale_color_manual(values = c("green", "yellow", "orange", "blue", "red", "darkgreen")) #Dimensionsreduktion als Input #für ein Clusterverfahren library(umap) res_umap <- umap(d=x[, -1], n_neighbors = 20, min_dist = 0.05, n_components= 3) dim(res_umap\$layout) library(mclust) mc <- Mclust(res_umap\$layout, G = 10:20)</pre>

Hilfsmittel/Tools für Dimensionsreduktion / Clustering				
M	Idee/Ansatz	Voraussetzungen/Vorgehen	R-Code [Default]	
Skalierung	<p>Skalierung ändert Ergebnis</p> <ul style="list-style-type: none"> Beispiel PCA: Bei grossen Unterschieden in absoluten Varianz wird unskalierten Daten Richtung v. ersten Hauptkomponenten durch Variable mit grosser Varianz bestimmt. Oftmals lohnt es sich eine skalierte Version gegenüber einer unskalierten Version durchzuführen (bspw. bei einer Dimensionsreduktion). Dabei ist auf die Unterschiede in den anschliessenden Resultate zu achten. Welchen Einfluss hat die Skalierung? 	<p>Faustregeln: Skalieren,</p> <ul style="list-style-type: none"> Variablen verschiedene Einheiten haben (z.B. kg, m) wenn explizit gewollt ist, dass alle Variablen das gleiche Gewicht haben. <p>Nicht skalieren,</p> <ul style="list-style-type: none"> wenn alle Variablen die gleiche Einheit haben und vergleichbar sind. Wenn man nicht weiss, was tun, ist es meistens besser standardisieren. 	<pre>scale(dat)</pre>	

Unüberwachtes Lernen: Clustering (Beschreibende Methode)

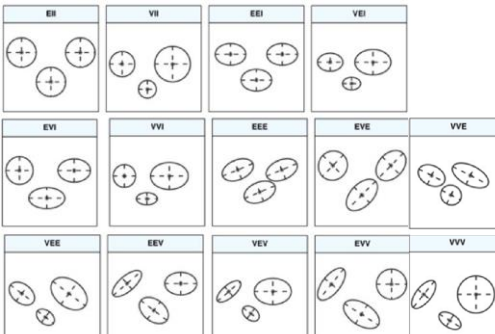
Ziel	Typische Schritte einer Clusteranalyse	Unterscheidung Clusterverfahren
<ul style="list-style-type: none"> • Set von Objekten in homogene Gruppen/Cluster unterteilen. Objekte im gleichen Cluster sollen untereinander ähnlicher sein als zu Objekten in anderen Clustern. • Um Clusteranalyse durchzuführen, muss ein Ähnlichkeits-/Distanzmass, für Objektevergleich, definiert werden. 	<ul style="list-style-type: none"> • Schritt 1: Explorative Datenanalyse und Datenbereinigung • Schritt 2: Skalierung Daten • Schritt 3: Wahl geeigneter Clusteralgorithmus (verschiedene Ansätze testen) • Schritt 4: Evaluierung geeigneten Anzahl Cluster • Schritt 5: Analyse der resultierenden Cluster 	<ul style="list-style-type: none"> • Partitionsverfahren (Raum unterteilen): Unterteilung der Objekte in sich nicht überschneidende Cluster, so dass jedes Objekt in genau einem Cluster ist. • Hierarchisches Clustern: Cluster werden innerhalb von grösseren Clustern unterteilt und bilden Baum von Cluster

	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
K-Means / WCV (within-cluster variation)	<ul style="list-style-type: none"> • Unterteilen der Beobachtungen in k homogene Cluster, so dass die totale Variation innerhalb der Cluster (WCV within-cluster variation) über alle K Cluster, so klein wie möglich ist. • $\min_{C_1, \dots, C_k} (\sum_{k=1}^K WCV(C_k))$ • WCV für einen Cluster k C_k, $WCV(C_k) = \frac{1}{\#C_k} \cdot \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$ • Verw. quadrierte euklidische Distanzen • $\#C_k$ Anz. Beobachtungen im Cluster k • p Anz. Features (Dimension) • Minimierung WCV nicht trivial • $\frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^N$ 	<ul style="list-style-type: none"> • Datenmatrix • Simpler guter Algorithmus (liefert lokale Lösung): <ol style="list-style-type: none"> 1. Beobachtungen zufälliger Gruppe zuordnen 2. Für jeden Cluster wird das Zentrum bestimmt. 3. Alle Beobachtungen werden Clusterzentrum zugeordnet, zudem sie kleinste Distanz haben. <ul style="list-style-type: none"> ○ Wiederholen Schritte 2-3 bis es keine Änderungen mehr gibt. ○ Lösung von Startkonfiguration abhängig. ○ Workaround: Algorithmus mehrmals durchführen und Simulation mit tiefstem WCV verwenden. • Algorithmus für div. k durchführen • WCV bestimmen und ihn vs. k auftragen. k nach letzten grossen Abfall bestimmen (Ellbogen-Methode) → Anzahl Cluster • Festlegen Anz. Clustern: Verschiedene Ansätze: Ellbogen-, Silhouetten-Methode, Domain knowledge, etc. → verschiedene Ansätze können zu versch. Anz. Cluster führen 	<p>+ Schnell und einfach</p> <ul style="list-style-type: none"> - Clusteranz. muss vorgegeben werden - Probleme mit verrauschte Daten/Ausreisser - Schwierigkeiten mit hochdim. Daten, führt zu wenig aussagekräftigen Ergebnissen - nicht fähig, Cluster mit nicht-konvexen Formen zu erkennen - kann in lok. Minima hängen bleiben - benötigt Datenmatrix - kann nur mit euklidischen Distanzen umgehen (andere nicht möglich). - Tendenz, «kugelförmige» Gruppen mit gleich vielen Objekten zu bilden. - Erkennt längliche Ausdehnung/verschiedene Streuungen in Gruppen schlecht - bei Ausreissern können komische Effekte auftreten, falls nicht für jeden Ausreisser (oder für Anhäufung von Ausreissern) ein separates Cluster gebildet wird. - Lösung von Startkonfiguration abhängig 	<pre>set.seed(17) #verschiedene seeds testen cl <- kmeans(dat, centers = 3) plot(x, col = cl\$cluster, las = 1) points(cl\$centers, ch = 2, cex = 2) #Zentrum cl\$withinss #WCV Wert pro Klasse #Visualisierung bei mehr als 2 Dimensionen: 2D-Plot mit PCA, MDS, t-SNE/UMAP # Bestimmung k k <- 20; wcv <- rep(0, k) for (i in 1:k) wcv[i] <- sum(kmeans(dat, centers = i)\$withinss) plot(1:k, wcv, type = "b", xlab = "Anzahl Cluster", ylab = "WCV") #Ellbogen-Plot # Alternative library(factoextra) fviz_nbclust(dat, kmeans, method = "wss", k.max = 10)</pre>
K-medoids / PAM	<ul style="list-style-type: none"> • PAM (Partitioning Around Medoids) • Sehr ähnlich zu k-means • Finden repräsentative Objekte, genannt Medoids, in Clustern • Ausgehend von einer Distanzmatrix wird über alle Objekte die Gesamtdistanz der Objekte zu den Objekten im gleichen Cluster minimiert $\min(\sum_{i=1}^n \sum_{i'=j}^n d(i, j) \cdot z_{ij})$ • WCV-Ansatz funktioniert nicht, da Kriterium nicht monoton abnehmend • Max. durchschnittliche Silhouettenbreite kann auch hier bestimmt werden 	<ul style="list-style-type: none"> • Distanzmatrix • Algorithmus <ol style="list-style-type: none"> 1. k Objekte werden zufällig als Clusterzentren gewählt 2. Jedes Objekt wird dem nächsten Clusterzentrum zugewiesen 3. Kostenfunktion wird evaluiert 4. Iterativ wird für jeden Cluster untersucht, ob Medoid Vertauschen 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>Vorteile gegenüber K-Means</p> <ul style="list-style-type: none"> + robuster gegenüber Ausreissern + Distanzmatrix + geht mit allen Distanzmassen + Medoids als repräsentative Objekte v. Cluster (hilft zu interpretieren) <ul style="list-style-type: none"> - Clusteranz. muss vorgegeben werden - Bestimmung über Gütekriterium (Ellbogen-Methode bei K-Means) nicht möglich, da dieses nicht monoton abnimmt (dafür kann dies aber mit Silhouettenbreite bestimmt werden) - langsam für grosse Daten (→ <code>clara()</code>) 	<pre>library(cluster); cl_pam <- pam(dist, k = 3) cmd <- cmdscale(dist, k = 2) #für Plot plot(cmd, col = cl_pam\$clustering) points(cl_pam\$medoids, pch = 8, cex = 2) cl_pam\$medoids #Medoid k <- 20; sil <- rep(0, k) for (i in 2:k){ sil[i] <- summary(silhouette(x = pam(dist, k = i)\$cluster, dist = dist))\$avg.width} plot(x = 1:k, sil, type = "b", xlab = "Anz Cluster", ylab = "mittlere Silhouettenbreite") abline(v = which.max(sil)) # Alternative (braucht Datenmatrix!): library(factoextra) fviz_nbclust(dat, pam, method = "silhouette") # Oder auch direkt mit library(fpc) cl_pamk <- pamk(dat, krange = 2:5)</pre>

<p>Silhouetten als Qualitätskontrolle</p> <ul style="list-style-type: none"> Silhouetten-Koeffizient einer Beobachtung i: $sil_i = \frac{b_i - a_i}{\max(a_i, b_i)} \in [-1, 1]$ a_i: mittlere Distanz zwischen Punkt i und allen anderen Punkten aus gleichen Cluster b_i: mittl. Distanz zwischen i und allen Punkten nächstgelegenen Nachbarcluster $b_i \gg a_i \rightarrow$ benachbarter Cluster weit weg $sil_i \rightarrow 1$ $b_i \approx a_i \rightarrow$ benachbarter Cluster nahe $sil_i \rightarrow 0$ \rightarrow Mittl. Silhouettenbreite kann auch Wahl optimalen Anzahl Cluster verwenden 	<p>Faustregel für mittlere Breite im Cluster</p> <ul style="list-style-type: none"> $0.70 < sil < 1.00$: Gute Struktur wurde gefunden $0.50 < sil < 0.70$: vernünftige Strukturen $0.25 < sil < 0.50$: Schwache Strukturen \rightarrow weitere Analyse $-1 < sil < 0.25$: Forget it! <p>Was bedeutet negativer Silhouettenkoeffizient?</p> <ul style="list-style-type: none"> Beobachtung ist im Mittel näher bei den Beobachtungen im benachbarten Cluster als zu denen im Eigenen. Clustering muss nicht zwingend schlecht sein. 		<p>R-Code</p> <p>K-Means</p> <pre>cl <- kmeans(dat, centers = 3) plot(silhouette(x = cl\$cluster, dist = dist(dat)))</pre> <p>PAM</p> <pre>cl_pam <- pam(dist, k = 3) plot(silhouette(x = cl_pam \$clustering, dist = dist))</pre>
---	--	--	---

	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
<p>Hierarchische Clusteranalyse</p>	<ul style="list-style-type: none"> Ein Set von verschachtelten Clustern, welche in einem hierarchischen Baum organisiert werden können. Anzahl Dendrogramme mit n Blättern: $(2n - 3)! / [(2(n - 2))(n - 2)!]$ Können nicht alle Möglichkeiten testen \rightarrow Brauchen Heuristik Top-Down (divisiv): <ul style="list-style-type: none"> Start: alle Objekte in einem Cluster Betrachte jede Möglichkeit die Daten in 2 Cluster zu unterteilen Wähle die Beste aus Repetieren für beide Cluster, usw. bis alle Objekte einen eigenen Cluster bilden. Bottom-Up (agglomerativ): <ul style="list-style-type: none"> Start: alle Objekte als eigene Cluster Finde das beste Paar zum verbinden Repetieren bis alle Cluster zu einem vereint sind \leftarrow übliche Form 	<p>Distanzmatrix</p> <p>Algorithmus, Start: Distanzmatrix</p> <ol style="list-style-type: none"> Jedes einzelne Objekt wird als Gruppe betrachtet. $\rightarrow n$ Gruppen Suche das kleinste Element in der Distanzmatrix. Entsprechende Elemente werden zu einer Gruppe verbunden und aus der Distanzmatrix entfernt. Berechne die Distanzen von der neuen Gruppe zu allen anderen. Wiederhole Schritte 2 und 3 bis nur noch eine Gruppe übrigbleibt. <p>Problem: Brauchen Verallgemeinerung, um Distanzen zwischen Clustern bestimmen. \rightarrow Linkage</p> <p>Distanzen zwischen Clustern: Linkage</p> <ul style="list-style-type: none"> Single-Linkage: kleinste Distanz zwischen den zwei Clustern (kleinste Distanz zwischen 2 Objekten) \rightarrow single Complete-Linkage: grösste Distanz zwischen 2 Objekten \rightarrow complete Average-Linkage: mittlere Distanz zwischen allen Objekten \rightarrow average Ward: versucht die Varianz der verbunden Cluster zu minimieren (hierarchische Version von K-Means) \rightarrow ward.D2 <p>Clusterergebnis hängt von der Datenstruktur, dem Distanzmass und der Linkage-Methode ab.</p> <p>Allgemein: Clustering ist exploratives Tool!</p> <p>\rightarrow Wählen Linkage Methode, welche «beste» Resultat produziert!</p>	<p>+ Verschiedene Distanzen möglich (euklidisch, Manhattan, Gower, ...)</p> <p>Nachteile</p> <p>-</p> <p>Dendrogramm</p> <ul style="list-style-type: none"> Position des Verbindungsknotens auf Abstandsskala gibt Abstand zwischen Clustern an (abhängig von Verknüpfungsmethode). Wenn zwei Cluster auf einer Höhe von bspw. 22 verbunden wurden, bedeutet das, dass der Abstand zwischen diesen Clustern 22 beträgt. Wollen feststellen, wo Abstand zwischen den Clustern, die kombiniert werden, gross ist, um Cluster zu unterteilen. Wir vergleichen die Abstände zwischen sequenziellen Verbindungsknoten (hier vertikale Linien). Achtung! Interpretation der Abstände zwischen den Clustern nicht zulässig. Um sinnvolle Cluster zu finden, wird das Dendrogramm bei einem Schwellwert geschnitten. Dazu betrachtet man die Differenzen der Cluster-Distanzen zwischen den Aufsplitterungen. Üblicherweise sind Differenzen zuunterst am kleinsten und oben oft gross. Dazwischen gibt es meistens eine Zone, wo sich Differenzen plötzlich stark verkleinern. Oberhalb davon ist geeigneter Ort, um den Schwellwert zu setzten. <p>Erhalt Schwellenwert:</p> <pre>hc\$height[length(hc\$height):1] diff(hc\$height)[(length(hc\$height)-1):1]</pre>	<pre>hc <- hclust(dist(dat), method = "ward.D2"); plot(hc, las = 1) #In k = 3 Cluster unterteilen/einzeichnen rect.hclust(hc, k = 3, border = "red") #Schwellenwert sieh +/- grhc <- cutree(hc, k = 3) #Einteilung Beobachtung in Cluster grhc <- cutree(hc, h = 5) #Einteilung mit Höhe des Schnittes # Plotten der Cluster mit MDS cmd <- cmdscale(dist); plot(cmd, xlab = "", ylab = "", col = c("green", "blue", "red")[cutree(hc, h = 8)]) text(cmd, labels = rownames(cmd), col = c("green", "blue", "red")[cutree(hc, h = 8)]) plot(silhouette(x = grhc, dist = dist)) k <- 20; sil <- rep(0, k) for (i in 2:k){ sil[i] <- summary(silhouette(x = cutree(hc, k = i), dist = dist))\$avg.width} plot(x = 1:k, sil, type = "b", xlab = "Anz Cluster", ylab = "mittlere Silhouettenbreite"); abline(v = which.max(sil)) #Alternative: Funktion agnes aus cluster</pre>
	<p>Heatmap</p> <ul style="list-style-type: none"> Visualisierung einer Datenmatrix blaue Werte \rightarrow kleine Werte rote Werte \rightarrow grosse Werte Permutation der Zeilen und Spalten, so dass ähnliche Objekte nahe sind Erlaubt zu verstehen, in welchen Merkmal sich Cluster unterscheiden Farbbalken erlauben das Überprüfen von Hypothese über Assoziationen zwischen Clustern und externen kategorialen Variablen, die beim Clustering nicht verwendet wurden. 			<pre># Heatmap dat_num <- dat[,3:8] x <- t(as.matrix(dat_num)) colnames(x) <- dat\$Car heatmap(x) # Default # distfun = dist -> euclidean, # hclustfun = hclust -> linkage = complete, # scale = "row" -> Zeilen standardisiert (mean = 0, sd = 1) # Elegante Alternative library(pheatmap) pheatmap(x, scale = "row", angle_col = 45) # Ergänzen mit Meta-Daten an <- data.frame(Country = dat\$Country) rownames(an) <- colnames(x) pheatmap(x, scale = "row", annotation_col = an)</pre>

Partitionierungsmethoden (K-Means, PAM-Clustering) und hierarchisches Clustering funktionieren gut bei kompakten und gut getrennten Clustern. Mit komplexeren Daten und insbesondere Rauschen und Ausreissern in den Daten können sie nur schlecht umgehen. Insbesondere Ausreisser sind für Daten typisch. \rightarrow Alternative DBSCAN - Density-Based Spatial Clustering of Applications with Noise

Modellbasiertes Clustering	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]																																																																											
	<ul style="list-style-type: none"> Idee: Daten stammen aus Mischverteilung, wobei Daten aus jedem Cluster aus anderen Verteilung stammen. Punkte können mit gewissen Wahrscheinlichkeit einem Cluster zugeordnet werden. Statistisches Modell für Beschreibung Clusterform. Populärste Annahme multivariate Normalverteilung für Cluster (für stetige Daten) Multivariate Normalverteilung (Verallgemeinerung der Normalverteilung auf mehrere Dimensionen) Bestimmt durch Verteilungsparameter: <ul style="list-style-type: none"> Erwartungswertvektor $\mu(\mu_{x1}, \mu_{x2}, \dots)$ Kovarianzmatrix Σ <p>Expectation-Maximization (EM-Algorithmus):</p> <ul style="list-style-type: none"> Gesucht n Cluster aus multivariaten Normaldichten mit $E(X) = \mu_j$ und Kovarianzmatrix Σ_j ($j = 1, \dots, n$) Clustergrößen in Anteilen $\pi_1, \pi_2, \dots, \pi_n$ (wobei $\sum_j(\pi_j) = 1$) Parameterschätzung (μ_j, Σ_j, π_j) aus Daten Maximierung Log-Likelihood mit iterativem EM-Algorithmus <p>Da Clustering auf Modell basiert, kann Unsicherheit für einzelnen Beobachtungen analysiert werden.</p>	<ul style="list-style-type: none"> Datenmatrix Start Algorithmus Zufällige Startparameter (μ_j, Σ_j, π_j) E-Step <ul style="list-style-type: none"> r_{ic}: W'keit, dass Punkt x_i zu Cluster c gehört Berechnung <i>Likelihood · Prior</i> für Model c und Normalisieren der Gewichte über alle Cluster. Je wahrscheinlicher x_i zum c-ten Cluster gehört, desto grösser r_{ic} M-Step <ul style="list-style-type: none"> Für jeden Cluster c werden Parameter anhand gewichteten Datenpunkte aktualisiert $m_c = \sum_i(r_{ic}), \quad \pi_c = \frac{m_c}{m}, \quad \mu = \frac{1}{m_c} \sum_i(r_{ic}x^{(i)})$ $\sigma_c = \frac{1}{m_c} \sum_i r_{ic} (x^i - \mu_c)^T (x^i - \mu_c)$ Rechenaufwand Kovarianzmatrix jedes Clusters hat Dim. $p \times p$. Falls p gross (oder n) müssen sehr viele Parameter geschätzt werden. Wird instabil → Vereinfachung durch Restriktion Clustermodelle einfachste Restriktion $\Sigma_j = \sigma^2 I$, für $j = 1, \dots, n$ Einheitsmatrix; σ^2 Varianz → alle Cluster sind sphärisch (kugelförmig) mit gleichem Radius Bestes Modell? Optimales Modell und optimale Anzahl Cluster kann über maximales BIC (Bayes'sches Informationskriterium) ermittelt werden. $BIC_{mclust} = 2 \cdot I_{M_c}(x_1, \dots, x_n \hat{\theta}_{mle}) - v \cdot \log(n)$ I_{M_c}: Log-Likelihood des Modells M_c mit Max-Likelihoodschätzung (MLE) geschätzten Param $\hat{\theta}_{mle} = (\hat{\mu}_{mle}, \hat{\Sigma}_{mle})$ v: Anz. zu schätzender Parameter n: Anzahl Datenpunkte Mass für Effizienz des Modells Modellkomplexität (Parameteranz.) wird bestraft 	<p>+ Da Clustering auf Modell basiert, kann Unsicherheit für einzelnen Beobachtungen analysiert werden.</p> <p>- Benötigt Datenmatrix</p> <p>Vor- und Nachteile</p> <p>Verschiedene Modelle</p> <table border="1" data-bbox="1249 359 1675 890"> <thead> <tr> <th>Modell</th> <th>Verteilung</th> <th>Grösse</th> <th>Form</th> <th>Orientierung</th> </tr> </thead> <tbody> <tr><td>EII</td><td>sphärisch</td><td>gleich</td><td>gleich</td><td>—</td></tr> <tr><td>VII</td><td>sphärisch</td><td>variable</td><td>gleich</td><td>—</td></tr> <tr><td>EEI</td><td>diagonal</td><td>gleich</td><td>gleich</td><td>Achsen</td></tr> <tr><td>VEI</td><td>diagonal</td><td>variable</td><td>gleich</td><td>Achsen</td></tr> <tr><td>EVI</td><td>diagonal</td><td>gleich</td><td>variable</td><td>Achsen</td></tr> <tr><td>VVI</td><td>diagonal</td><td>variable</td><td>variable</td><td>Achsen</td></tr> <tr><td>EEE</td><td>elliptisch</td><td>gleich</td><td>gleich</td><td>gleich</td></tr> <tr><td>EVE</td><td>elliptisch</td><td>gleich</td><td>variable</td><td>gleich</td></tr> <tr><td>VEE</td><td>elliptisch</td><td>variable</td><td>gleich</td><td>gleich</td></tr> <tr><td>VVE</td><td>elliptisch</td><td>variable</td><td>variable</td><td>gleich</td></tr> <tr><td>EEV</td><td>elliptisch</td><td>gleich</td><td>gleich</td><td>variable</td></tr> <tr><td>VEV</td><td>elliptisch</td><td>variable</td><td>gleich</td><td>variable</td></tr> <tr><td>EVV</td><td>elliptisch</td><td>gleich</td><td>variable</td><td>variable</td></tr> <tr><td>VVV</td><td>elliptisch</td><td>variable</td><td>variable</td><td>variable</td></tr> </tbody> </table> 	Modell	Verteilung	Grösse	Form	Orientierung	EII	sphärisch	gleich	gleich	—	VII	sphärisch	variable	gleich	—	EEI	diagonal	gleich	gleich	Achsen	VEI	diagonal	variable	gleich	Achsen	EVI	diagonal	gleich	variable	Achsen	VVI	diagonal	variable	variable	Achsen	EEE	elliptisch	gleich	gleich	gleich	EVE	elliptisch	gleich	variable	gleich	VEE	elliptisch	variable	gleich	gleich	VVE	elliptisch	variable	variable	gleich	EEV	elliptisch	gleich	gleich	variable	VEV	elliptisch	variable	gleich	variable	EVV	elliptisch	gleich	variable	variable	VVV	elliptisch	variable	variable	variable	<pre>library(mclust) # spezifisches Modell mc <- Mclust(dat, modelNames = "EII") mc <- Mclust(dat) # alle Modelle mc\$modelName # bestes Modell # Visualisierung BIC aller Modelle plot(mc, what = "BIC") # Cluster plot(mc, what = "classification") # Alternative Visualisierung library(factoextra) fviz_mclust(mc, "BIC", palette = "jco") fviz_mclust(mc, "classification", geom = "point", palette = "jco") # Unsicherheit mc\$uncertainty plot(mc, what = "uncertainty") fviz_mclust(mc, "uncertainty", palette = "jco", xlab = "x1", ylab = "x2") # → grössere Symbole zeigen unsicherere Beobachtungen</pre>
Modell	Verteilung	Grösse	Form	Orientierung																																																																											
EII	sphärisch	gleich	gleich	—																																																																											
VII	sphärisch	variable	gleich	—																																																																											
EEI	diagonal	gleich	gleich	Achsen																																																																											
VEI	diagonal	variable	gleich	Achsen																																																																											
EVI	diagonal	gleich	variable	Achsen																																																																											
VVI	diagonal	variable	variable	Achsen																																																																											
EEE	elliptisch	gleich	gleich	gleich																																																																											
EVE	elliptisch	gleich	variable	gleich																																																																											
VEE	elliptisch	variable	gleich	gleich																																																																											
VVE	elliptisch	variable	variable	gleich																																																																											
EEV	elliptisch	gleich	gleich	variable																																																																											
VEV	elliptisch	variable	gleich	variable																																																																											
EVV	elliptisch	gleich	variable	variable																																																																											
VVV	elliptisch	variable	variable	variable																																																																											

	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
Dichtebasiertes Clustering DBSCAN	<ul style="list-style-type: none"> Idee: Erkennen von Gebieten mit hoher Datenpunktdichte als Cluster. Dichte Datenpunkts wird durch Anzahl der Punkte in seiner Umgebung ϵ geschätzt. ϵ-Umgebung besteht aus allen Punkten, die höchstens den Abstand ϵ haben → Parameter ϵ Mit Parameter <code>minPts</code> wird festgelegt, wie viele Punkte in der ϵ-Umgebung von einem Punkt p liegen müssen (Zentrumspunkt wird dazu gezählt), damit p zu einem Cluster gehört. <p>Drei Kategorien von Punkten:</p> <ul style="list-style-type: none"> Kernpunkt: Anzahl Datenpunkte in ϵ-Umgebung beträgt mindestens <code>minPts</code> (inklusive Kernpunkt). Randpunkt: Kein Kernpunkt, liegt aber in ϵ-Umgebung eines Kernpunktes. Rauschpunkt: weder Kern-/Randpunkt. → Cluster: Verbundene Kernpunkte plus dazugehörige Randpunkt. <p>Bestimmung Parameter</p> <ul style="list-style-type: none"> minPts: min. 3, je grösser Datensatz, desto grösser. Grössere Werte sind bei Datensätzen mit Rauschen und vielen Duplikaten meist besser und ergeben mehr signifikante Cluster. Wie gross sollten die Cluster mindestens sein? eps: kann mithilfe k-Nächste-Nachbarn-Graph abgeschätzt werden. Werden mittleren Abstände aller Punkte zu seinen k nächsten Nachbarn berechnet (mit $k = \text{minPts} - 1$). Abstände aufsteigender Reihenfolge auftragen. ϵ-Wert dort, wo Plot Ellbogen zeigt. Wenn ϵ zu klein, wird grosser Teil nicht geclustert. Zu hohem Wert verschmelzen Cluster. Kleine ϵ Werte sind vorzuziehen → Domainknowledge ist für Parameterbestimmung nützlich. 	<ul style="list-style-type: none"> Distanzmatrix <p>DBSCAN (Density-Based Spatial Clustering of Applications with Noise) Algorithmus:</p> <ol style="list-style-type: none"> Distanzmatrix über alle Punkte. Zufälliger (noch nicht besuchter) Punkt auswählen. Benachbarte Punkte mit Distanz $\leq \epsilon$ identifizieren. Weniger als <code>minPts</code> Punkte in der Nachbarschaft → Rauschpunkt (kann sich später wieder ändern) ⇒ 2. Mindestens <code>minPts</code> Punkte → Cluster wird gestartet. Punkt → Kernpunkt; bildet mit Punkten in der Nachbarschaft Cluster. Restliche Punkte im Cluster → Randpunkte. Iteration über alle Randpunkte. Falls in deren Nachbarschaft ebenfalls min. <code>minPts</code> vorhanden → Kernpunkt und Cluster wird erweitert. Prozess bis alle Punkte des Clusters gefunden wurden. ⇒ 2 Ende Algorithmus, wenn alle Punkte besucht. Benachbarte Kernpunkte mit ihren Randpunkten bilden Cluster. 	<ul style="list-style-type: none"> + Keine Vorgabe von Clusteranzahl + Filtert Rauschen i.d.R. gut aus Daten + Separate Ausweisung der Ausreisser + Beliebige Form/Grösser der Cluster + Benötigt nur Distanzmatrix <ul style="list-style-type: none"> - Schwierigkeiten mit Erkennung von Clustern unterschiedlicher Dichten. - Probleme bei hochdimensionalen Daten, da hier Dichte schwierig zu definieren ist (Curse of dimensionality). - Bestimmung von ϵ teilweise schwierig; erfordert Domänenwissen. - Parameter müssen bestimmt werden 	<pre>library(dbscan) db <- dbscan(dist, eps = 0.15, MinPts = 5) # Output, Cluster Nr. 0 = Ausreiser # 0 1 2 3 4 5 # 31 410 405 104 99 51 #Mit Dimreduktion 2-D Plot erstellen: plot(df\$x, df\$y, col = db\$cluster + 1) # Bestimmung Parameter library(dbscan) kNNdistplot(x = dist, k = 3)</pre>

Dilemma beim Clustern	Analyse der Cluster: Qualitätscheck
<ul style="list-style-type: none"> • Was ist ein Cluster? • Welche Feature/Merkmale/Variablen sollen verwendet werden? • Sollen die Daten normalisiert werden? • Gibt es Ausreisser in den Daten? Was haben die für einen Effekt? • Wie definieren wir die paarweise Ähnlichkeit? • Wie viele Cluster gibt es in den Daten? • Welche Clustering-Methode sollte verwendet werden? • Sind die entdeckten Cluster valid? <p>➔ Alle diese Fragen haben normalerweise keine klare Antwort!</p>	<ul style="list-style-type: none"> • Visualisierungen: <ul style="list-style-type: none"> ○ Dendogramme ○ Heatmap (für numerische Daten; Interpretation der Farben) ○ farbliche Clustervisualisierung in 2D Plot (PCA, MDS, t-SNE, UMAP) • Betrachte Zentrum oder repräsentatives Element (PAM) der Cluster • Silhouette-Plot • Quantifizierung der WCV within-cluster variation • Konfusionsmatrix (benötigt Labels / Ground truth Informationen; Clustern dann weniger interessant) • Analyse Featuresrelevanz → Klassifikation mit Cluster-ID als Label und betrachte Variablen Importance (vgl. Klassifikation) • Domainwissen beziehen! → generell schwierig in hohen Dimensionen

Überwachtes Lernen: Allgemeines zur Klassifikation

Ziel	Idee	Feature extrahieren	Ausgangslage	Ablauf	Fallstricke
Für Objekte, deren Klassenzugehörigkeit unbekannt ist, aufgrund der Merkmale (Variablen, Features) die richtige Klasse vorherzusagen.	<ul style="list-style-type: none"> • Trainiere Klassifizierer anhand von Trainingsdaten, bei welchen Klassen bekannt sind. • Verwende Klassifizierer um neue Beobachtung ohne Label zu klassifizieren. • Haben Trainingsdaten und unbekannte Daten respektive Testdaten. • Training/Validierung: Verwenden, um Modell trainieren und Hyperparameter bestimmen. • Testdaten: Zur Evaluation der Performance des finalen Modells 	<ul style="list-style-type: none"> • Wollen geeignete Feature/Merkmale aus Daten extrahieren, was nicht immer einfach ist. • Verschiedene Merkmale können kombiniert werden, um neue Merkmale zu generieren (z.B. PCA, UMAP Embedding). • Feature Engineering 	<ul style="list-style-type: none"> • Haben X und Y gegeben. Y Labels sind kategoriell und X ist Datenmatrix mit diversen Variablen. • In Klassifikation versuchen wir Labels Y anhand Features X vorherzusagen. • In Regression ist Zielvariable Y numerisch. 	<ol style="list-style-type: none"> 1. Zufälliges Aufteilen in Training- und Test-Daten (i.d.R. 80 % und 20 %) 2. Anpassen Modell auf Trainingsdaten (inkl. Preprocessing & Modellselektion) 3. Modellevaluation: Wie gut wird Modell bei neuen, nicht-gesehenen Daten performen? 	<ul style="list-style-type: none"> • Kein Trainingsverfahren (Datenaufbereitung gehört dazu) darf Testdaten sehen. Sperrn diese bis zum letzten Test weg! Grund: Wenn verschiedene Modelle anhand Testdaten evaluiert, besteht Gefahr, dass man zu optimistisch ist. Rückschlüsse der Modellperformance auf neue Daten können nicht gezogen werden. • Unsymmetrische Daten: Unausgeglichene Datensätze können Ergebnis verzerren (Bsp. binäre Variable, 1 → 90 %, 2 → 10 %)

Aufteilung Datensatz: `library(caret); index <- createDataPartition(dat$y, p = 0.8, list = F) #80 % Trainingsdaten, Output Index von Trainingsdaten` `trainDat <- dat[index,]` `testDat <- dat[-index,]`

Evaluationsmetrik	Berechnung	Was wird gemessen	Fehlerarten	Modellselektion
Accuracy/Genauigkeit	$\frac{TP+TN}{n_{test}}$	Standardmetrik, unterschiedliche Kosten für Fehler werden nicht berücksichtigt	<ul style="list-style-type: none"> • Trainingsfehler, naiver Fehler bzw. In-Sample-Fehler: Fehler an Trainingsdaten evaluiert • Verallgemeinerungsfehler, Testfehler bzw. Out-of-Sample-Fehler: Fehler an neuen, bisher ungesehenen Daten (aus der Stichprobe) evaluiert • Komplexeres Modell performt bei Trainingsdaten immer besser als einfaches Modell • Modelle sollten anhand von Out-of-Sample-Daten ausgewertet werden, um Verallgemeinerungsfehler zu ermitteln. • Beim Vergleich der Performance an Trainingsdaten sollte Modellkomplexität mitberücksichtigt werden (Strafe für Komplexität). • Bei zwei Modellen mit ähnlichen Verallgemeinerungsfehlern sollte das einfachere Modell dem komplexeren Modell vorziehen. • Overfitting Phänomen (Überanpassung): Modell passt gut zu Trainingsdaten (kleiner Trainingsfehler/grosse Genauigkeit), hat aber hohen Testfehler/kleine Genauigkeit. 	<ul style="list-style-type: none"> • Welches Modell sollten wir wählen? • Vergleich unterschiedlicher Modelle. • Oft gibt es einen «Hebel» (Hyperparameter), der die Komplexität/Performance eines Klassifikators steuert. • Welche Merkmale sollten verwendet werden? • Soll ein Merkmal transformiert werden (Wurzel, Logarithmus, etc.)?
Precision/Präzision Pos Pred Value	$\frac{TP}{TP+FP}$	Wie viele vorhergesagten ⊕-Objekte sind richtig? → wollen hohen TP, tiefen FP		
Recall/Sensitivität True-Positive-Rate	$\frac{TP}{TP+FN}$	Wie viele ⊕-Objekte wurden richtig vorhergesagt? → wollen hohen TP, tiefen FN		
Spezifität True-Negative-Rate	$\frac{TN}{TN+FP}$	Wie viele der ⊖-Objekte wurden richtig vorhergesagt?		
F1-Score	$\frac{2 \cdot recall \cdot precision}{recall + precision}$	Harmonisches Mittel v. Recall % Precision		
Kappa	$\frac{accuracy - E[accuracy]}{1 - E[accuracy]}$	Relative Verbesserung im Vergleich zu einem zufälligen Prädiktor		
$E[accuracy] = \frac{TP+FP}{n_{test}} \cdot \frac{TP+FN}{n_{test}} + \frac{TN+FN}{n_{test}} \cdot \frac{TN+FP}{n_{test}}$, $n_{test} = TP + FP + FN + TN$				
Evaluationsmetriken: Je nach Fragestellung kann eine andere Metrik wichtig sein.				

Mod.validierung	Einfache Validierung	Leave-One-Out KV LOOKV	K-fachen KV – KKV	Bootstrap-Valid.	Out of Bag	R-Code Trainieren [Def]
<ul style="list-style-type: none"> • Wie abschätzen, wie gut Modell mit ungesehenen Daten ist? <ul style="list-style-type: none"> ○ Wichtig: Ohne Verwendung Testdaten. ○ Modellvergleiche auf Testdatenbasis führen zu überoptimistischen Fehler-schätzungen! → Validierung bereits in Trainingsdaten erforderlich → Welchen Ansatz? Verwenden (mehrere) 10-fache KV/LOOKV 	<ul style="list-style-type: none"> • Aufteilung Gruppen: Trainings- und Validierung <ul style="list-style-type: none"> – Schätzfehler sind variabel abhängig, welche Beobachtungen in Trainings-/Testmenge sind. – Nur Teilmenge Beobachtungen (Training) wird zur Modellanpassung verwendet. Statistische Methoden schneiden schlechter ab, wenn sie auf weniger Beobachtungen trainiert werden. – Testfehler für Modellanpassung an gesamten Datensatz wird daher oft überschätzt. 	<ol style="list-style-type: none"> 1. Teilen Daten in n Datensätze auf, wobei Trainingsdaten aus allen Punkten ausser einem bestehen. 2. Modell wird einmal mit allen n Datensätzen trainiert und dann am ausgeschlossenen Punkt validiert. 3. Schätzen Testfehler aus Anzahl Fehlklassifizierungen z.B. Anteil Fehlklassierung: $CV_{(n)} = \frac{1}{n} \sum_{i=1}^n Err_i$, $Err_i = I(y_i \neq \hat{y}_i)$ <ul style="list-style-type: none"> + Nicht zufällig in Aufteilung zwischen Trainings- und Validierungsdaten. + Kleinster Bias bei Schätzung Testfehler (Verwendung fast aller Daten zur Schätzung) – Numerisch komplex/langsam (Modell n-mal ausgeführt) 	<ol style="list-style-type: none"> 1. Unterteilen Daten in k sich nicht überschneidende Gruppen. 2. Modell k-mal trainiert für alle Gruppen ausser einer, und jedes Mal mit dieser Gruppe validiert. 3. Schätze Testfehler aus Anz. falsch klassifiziert Beobachtung: $CV_{(k)} = \frac{1}{k} \sum_{i=1}^k Err_k$ <ul style="list-style-type: none"> + Schneller als LOOKV – Ergebnis weist gewisse Zufälligkeit auf (abhängig Datenaufteilung in Gruppen) –/+ Etwas grösserer BIAS als LOOKV. I.d.R. unproblematisch und besser, da geringere Varianz als LOOKV. → 10-fache KV Standard 	<ul style="list-style-type: none"> • Erzeugt Stichproben aus vorhand. Daten. • Gleiche Grösse n wie ursprüngliche Stichprobe, ziehen mit Zurücklegen $x = (x_1, \dots, X_n)$ • Benutze Datensatz als Ersatz für unbekannte zugrundeliegende Verteilung und ziehe daraus Stichproben. • Stichprobe aus Trainingsdaten. • Validierung auf Daten, die nicht in Bootstrapstichprobe 	<p>Evaluation Random Forest</p> <ul style="list-style-type: none"> • oob Fehlerrate: Bei jedem Baum werden nicht verwendeten Testdaten jeweils zur Evaluation der Fehlerrate verwendet. • Mitteln über alle Bäume → KV on the fly! 	<p>Für Steuerung der KV <code>trainControl()</code> <code>ctrl <- trainControl(method = ...)</code></p> <ul style="list-style-type: none"> • "LOOCV": LOOKV (<code>Parameter</code>) • "cv": reg. k-KV (<code>number = k</code>) [10] • "repeatedcv": rep k-KV (<code>repeats = Anz. Rep. KVs, number = k</code>) • ["boot"]: (<code>number = Anz Rep</code>) [25] • "oob": out of Bag RandForest • "none": nichts <p><code>train(x, y, method = "knn", trControl = ctrl, tuneGrid = data.frame(k = 1:10), preProcess = c("center", "scale"))</code></p> <p><code>tuneGrid</code>: Testparameter (alternativ <code>tuneLength</code>)</p> <p><code>center</code>: Daten zentrieren</p> <p><code>scale</code>: Skalierung mit <code>sd()</code></p> <p><code>pca</code>: PCA mit <code>thresh</code>: Anteil Varianz, <code>pcaComp</code> Anz. PCs</p>

Überwachtes Lernen: Klassifizierer

No-Free-Lunch-Theorem				
<ul style="list-style-type: none"> • Modell ist vereinfachte Version der Realität: Vereinfachung, um überflüssige Details eliminieren und sich auf Aspekte zu konzentrieren, die man verstehen will. • Vereinfachungen beruhen auf Annahmen; Annahmen können in einigen Situationen zutreffen, in anderen nicht. Modell, das bestimmte Situation gut erklärt, kann in anderen Situation versagen. • Theorem besagt, dass es kein Modell gibt, das für jedes Problem am besten funktioniert. Annahmen eines Modells für ein Problem gelten möglicherweise nicht für andere Probleme. Beim maschinellen Lernen ist es üblich, mehrere Modelle auszuprobieren, um das zu finden, das am besten funktioniert. ⇒ Darum verschieden Ansätze (Klassifizierer) 				
Idee/Ansatz	Voraussetzungen/Vorgehen	Pro (+) / Contra (-)	R-Code [Default]	
k-Nächste-Nachbarn KNN	<ul style="list-style-type: none"> • Ausgangslage: Objekt x_0 mit Klassenlabel unbekannt. • Finde die k Trainingsobjekte mit der kleinsten Distanz zu x_0 (eukld. Dist. → Daten ggf. skalieren!) • Verwende die Klasse der Mehrheit der k Nachbarn als Klassenlabel für x_0 <p>Was ist die richtige Modellkomplexität?</p> <ul style="list-style-type: none"> • Wollen Überanpassung vermeiden, aber dennoch ausreichend komplex sein. • Ein kleines k führt zu einem komplexen Modell, ein grosses k bedeutet ein einfaches Modell. • Bei zwei Klassen sollte immer ungerade Anzahl nächste Nachbarn wählen, damit eindeutige Entscheidung getroffen wird. Sonst wird in Pattsituation Klassenzugehörigkeit ausgelost. Bei > 2 Klassen ist es nicht immer möglich, Pattsituation mit $k > 1$ zu verhindern 	<ul style="list-style-type: none"> • Skalierung der Daten • Beim knn-Klassifizierer, so wie bei vielen anderen ML-Ansätzen, spielt es eine Rolle, welche Grössenordnung die einzelnen Features habe. • Feature mit grösseren Skala haben mehr Einfluss. • Damit alle Feature den gleichen Einfluss haben, müssen die Daten skaliert werden. <p>Achtung: Bei $k = 1$ muss die Fehlerrate bei den Trainingsdaten immer 0 sein, da nur das eigene Label die Klasse definiert.</p>	<p>+</p> <ul style="list-style-type: none"> - Ist nur für euklidische Distanzen ausgelegt - Ist für Datensätze mit vielen Merkmalen weniger geeignet 	<pre>res.knn <- train(y = trainDat\$class, x = trainDat[, 1:4], method = "knn", trControl = ctrl, tuneGrid = data.frame(k = seq(1,15,2)) #tuneGrid = Welche k ausprobieren? plot(res.knn) #Darstellung Accuracy nach # Nachbarn confusionMatrix(predict(res.knn, testDat), reference = testDat\$class) #Alternative mit Formelschreibweise (FSW) res.knn <- train(form = class ~., data = trainDat, method = "knn") #FSW führt One-Hot-Codierung durch! #Alternative ohne caret mit knn Funktion: library(class) knn(train = trainDat[,-"label"], cl = trainDat\$label, test = newDat, k = k)</pre>
Naive Bayes-Klassifikator	<p>Bayes Klassifizierer:</p> <ul style="list-style-type: none"> • Verwendung insbesondere bei Textdaten (NLP) • Optimaler Klassifizierer: Beobachtung wird zur wahrscheinlichsten Klasse, gegeben Features, zugewiesen. ⇒ $P(Y = C_i X) \rightarrow$ bedingt Wahr'keit • Optimal falls Verteilung bekannt. • Bei welchen x-Wert liegt Entscheidungsgrenze? → Dort wo Wahrscheinlichkeit für beide Klassen gleich gross ist. ⇒ Perfekte Klassifikation nicht möglich <p>Wie sieht Entscheidungsgrenze aus? ($p = \# \text{ Var.}$)</p> <ul style="list-style-type: none"> • $p = 1 \rightarrow$ Entscheidungsgrenze: ein Punkt • $p = 2 \rightarrow$ Entscheidungsgrenze: eine Linie • $p = 3 \rightarrow$ Entscheidungsgrenze: eine 2D Ebene • $p = L \rightarrow$ Entscheidungsgrenze: L-1 Hyperebene • Likelihood kann im multivariaten Fall komplex sein. Braucht Vereinfachungen. → Naive Bayes <ul style="list-style-type: none"> • Problem wenn Häufigkeiten 0 (d.h. eine Kategorie einer Klasse kommt in Trainingsdaten nicht vor). → a posteriori wird 0 → Laplace Korrektur/Smoothing: konstanter Faktor (i.d.R. zwischen 1-2) wird bei Zählern dazu addiert, im Nenner wird Faktor mal Anz. Features dazu addiert. 	<ul style="list-style-type: none"> • Hauptsächlich für kategorielle Features Annahmen: <ul style="list-style-type: none"> • Alle Features sind voneinander unabhängig • Alle Features haben gleiches Gewicht. • Schätzung Likelihoods für Features über Häufigkeitsverteilung der Klassen in Trainingsdaten Theorem von Bayes für Klassifikation $P(Y = C_i X) = \frac{P(C_i) \cdot P(X Y=C_i)}{P(X)}$ <ul style="list-style-type: none"> • $P(C_i)$; a priori-Wahrscheinlichkeit (Anfangswahrscheinlichkeit) für Klasse i. → Bevor Daten sehen, klassifizieren gemäss dieser Wahrscheinlichkeit. • $P(Y = C_i X)$; beobachten Merkmalswert x und bestimmen damit a posteriori-Wahr'keit. • $P(X Y = C_i)$ Likelihood; Wahrscheinlichkeit für Beobachtung der Daten gegeben die Klasse. • $P(X)$; Randwahrscheinlichkeit für X, Wahrscheinlichkeit für die Prädiktorvariablen. ⇒ Klassifizierung gem. höchster a posteriori-Wahr'keit Umgang mit stetigen Variablen <ul style="list-style-type: none"> • Können diskretisiert werden. • Je nach Implementierung direkt verwenden: <ul style="list-style-type: none"> ○ Annahme normalverteilt (Schätzung Parameter mit mean und var) → Gaussian Naive Bayes ○ Nichtparametrische Dichteschätzung 	<p>+ Hauptsächlich für kategorielle Features (x_1, \dots, x_j) ausgelegt</p> <ul style="list-style-type: none"> - Annahme Unabhängigkeit sehr stark, aber Ansatz funktioniert häufig überraschend gut. 	<pre>ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3) nb <- train(y = trainDat\$Species, x = trainDat[, 1:4], method = "naive_bayes", trControl = ctrl, tuneGrid = data.frame(laplace = 0, usekernel = c(FALSE, rep(TRUE, 3)), adjust = c(0, seq(1, 2, 0.5)))) # tuneGrid, data.frame() ermöglicht es unterschiedlichen Features mit unterschiedlichen Schätzern laufen lassen # laplace: Faktor für Laplace-Korrektur, 0 = keine Korrektur, ansonsten Wert zwischen 1-2 # usekernel: TRUE = empirische Dichteschätzung für numerische Daten, FALSE = Gaussian Naive Bayes # adjust: Bandbreite für Dichteschätzung, wenn usekernel = T, dann Bandbreite angeben, bei usekernel = F, dann adjust = 0 confusionMatrix(predict(nb, testDat), reference = testDat\$Species) predict(nb, newdata = data.frame(alter = "21-30", einkommen = "mittel", geschlecht = "Frau")) #Voraussetzung aufgrund bestimmter Features # Alternative mit library(e1071): naiveBayes(einkauf ~ ., data = einkauf[,1:3])</pre>

Überwachtes Lernen: Klassifizierer

	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
	Preprocessing für praktische Anwendung von LDA und QDA • Normalverteilung: Für Feature wird Normalverteilung vorausgesetzt. Kann helfen, univariate Verteilungen jeder Variable zu überprüfen und zu transformieren (z.B. log/Wurzel für Exp.verteilung) • Ausreisser entfernen: Ausreisser können die Schätzungen verzerren. Ggf. hilft es, Ausreisser zu entfernen.			
Lineare Diskriminanzanalyse LDA	<ul style="list-style-type: none"> • Verwendung Bayes Theorem • Verteilungen der Klassen ist normalverteilt • Ausgelegt auf numerische Features. • Setzt keine Unabhängigkeit voraus. • Zusammenhang wird modelliert. Parameterschätzung mit Trainingsdaten: <ul style="list-style-type: none"> • Erw.wert μ_I wird mit Mittelwert aller Trainingsdaten aus Klasse I geschätzt. $\hat{\mu}_I = \frac{1}{n_I} \sum_{i: y_i=I} x_i$ • Var. σ^2 wird als gewichteter Mittelwert Var. aller k-Klassen geschätzt. $\hat{\sigma} = \frac{1}{n-k} \sum_{I=1}^k \sum_{i: y_i=I} (x_i - \hat{\mu}_I)^2$ • π_I wird als Anteil der Trainingsbeobachtungen geschätzt, die der I-ten Klasse angehören Multivariate Merkmalsverteilung <ul style="list-style-type: none"> • jede Klasse wird durch multivariate Normalverteilung mit Erwartungswertvektor μ ($\mu_{x_1}, \mu_{x_2}, \dots, \mu_{x_p}$) und • Kovarianzmatrix Σ (alle Klassen gleich) definiert. 	<ul style="list-style-type: none"> • Numerische Features • Keine Parameter zum Tunen! • Gleiche Varianz: Es ist fast immer gute Idee, Daten vor Verwendung zu standardisieren. Vorgehen: 1. Schätzen Prior der Klassen gem. Klassenhäufigkeit im Trainingsdatensatz oder Vorwissen 2. Schätzen Parameter für zugrundeliegende Normalverteilungen (μ_1, μ_2, σ^2) → Annahme LDA: Prädiktoren für jede Klasse sind normalverteilt mit unterschiedlichem Erwartungswert, aber gleicher Varianz 3. Bestimmen Entscheidungsgrenze nach Bayes-Regel: $P(Y = C_I X) = \frac{\pi_I f_I(x)}{\sum_{i=1}^K \pi_i f_i(x)}$ mit $f_i(x)$: Dichte von X für Klasse k, π_I : Prior für Klasse I	+ Keine Unabhängigkeitsannahme + Methode funktioniert erstaunlich gut, selbst bei nichtnormalverteilten Gruppen. – Annahme Normalverteilung – Nur für numerische Features – Annahme: Varianzen aller Klassen sind gleich – Ist für Datensätze mit vielen Merkmalen weniger geeignet	<pre>lda_fit <- train(y= trainDat\$Species, x = trainDat[, 1:4], method = "lda", trControl = ctrl, preProcess = c("center", "scale"))</pre> <p># Vorhersage <code>confusionMatrix(predict(lda_fit, testDat), reference = testDat\$Species)</code></p> <p># Teilweise muss für DA zuerst Dimensionsreduktion durchgeführt werden: <code>res.pca <- prcomp(trainDat[, -1], scale = F)</code> <code>v_80 <- which(cumsum(res.pca\$sdev^2)/sum(res.pca\$sdev^2)>0.8)[1]</code></p> <pre>train_reduced <- data.frame(label = trainDat\$label, res.pca\$x[, 1:v_80]) m.llda <- train(label ~ ., data = train_reduced, method = "lda", trControl = ctrl)</pre>
Quadratische DA	<ul style="list-style-type: none"> • Erweiterung LDA auf QDA • Gibt für jede Klasse eigene Kovarianzmatrix Σ • X geht nun quadratisch in die Diskriminanzfunktion $\delta k(x)$ ein → Klassengrenzen können jetzt eine quadratische Form haben. • Bei QDA müssen wesentlich mehr Parameter geschätzt werden. 	<ul style="list-style-type: none"> • Numerische Features • Keine Parameter zum Tunen! Vorgehen: 1. Modellierung Verteilung Features X der versch. Gruppen individuell in Trainingsdaten 2. Bayes Theorem, um bedingte Wahrscheinlichkeit → $P(Y = C_I X)$ zu bestimmen. 3. Ordne die Beobachtung der Klasse mit der höchsten A-posteriori Wahrscheinlichkeit zu.	+ Keine Unabhängigkeitsannahme + Verschiedene Varianzen in Klassen möglich + Methode funktioniert gut, selbst bei nichtnormalverteilten Gruppen. – Annahme Normalverteilung – Nur für numerische Features – Hohe Berechnungskomplexität – Ist für Datensätze mit vielen Merkmalen weniger geeignet	<pre>qda_fit <- train(y = trainDat\$Species, x = trainDat[, 1:4], method = "qda", trControl = ctrl, preProcess = c("center", "scale"))</pre> <p># Vorhersage <code>confusionMatrix(predict(qda_fit, testDat), reference = testDat\$Species)</code></p>
Log. Regression	<ul style="list-style-type: none"> • $\log\left(\frac{p(y x)}{1-p(y x)}\right) = \beta_0 + \beta_1 \cdot x$ • Schätzung Parameter $\hat{\beta}_0, \hat{\beta}_1$ mit Max. Likelihood $I: I(\hat{\beta}_0, \hat{\beta}_1) = \prod_{i: y_i=1} p(y_i x_i) \prod_{i: y_i=0} (1 - p(y_i x_i))$ • Linearer Prädiktor (rechte Seite) kann weitere erklärende Variablen haben (vgl. multiple Regression) • Logistische Regression hat keinen Fehlerterm. Streuung Ergebnisse wird Bernoulli-Verteilung modelliert. • Geschätzten Modellparameter kann man analysieren. • Koeffizienten sind log-Odds. Können diese in Odds umwandeln. Z.B. <code>exp(res.glm\$coefficients[2])</code> 	<ul style="list-style-type: none"> • kein Tuningparameter Multinomiale log Regression (Mehr als 2 Klassen) <ul style="list-style-type: none"> • Eine der K Klassen wird zur «Referenzkategorie» k_0 (erste Faktorstufe, optimal wenn es grösste Klasse ist) • Für alle anderen $K - 1$ Klassen wird eine separate logistische Regression gegenüber der Referenzkategorie modelliert. $\log\left(\frac{P(Y_i=k X_i)}{P(Y_i=0 X_i)}\right) = \text{logit}(p_k)$ • Gilt $\sum_{k=0}^{K-1} p_k = 1 \rightarrow p(Y_i = k X_i) = \frac{\exp(\beta_k \cdot X_i)}{1 + \sum_{j=1}^{K-1} \exp(\beta_j \cdot X_i)}$ 	+ Methode geht auch für unbalancierte Daten, Konfusionsmatrix ist dann jedoch heikel + Gibt auch Ansätze für ordinale Variablen (<code>method = "polr"</code> in <code>caret</code>). + – Ist für Datensätze mit vielen Merkmalen weniger geeignet	<pre>res.glm <- train(y = trainDat\$y, x = trainDat\$x, method = "glm", family = "binomial")</pre> <p><code>summary(res.glm)\$coefficients</code></p> <pre>confusionMatrix(predict(res.glm, newdata = testDat), testDat\$Aktiv)\$table</pre> <pre>res.glm <- train(y = y, x = x, method = "multinom", trControl = ctrl, tuneLength = 7, trace = F)</pre>
Logistische Regression, k-NN und Diskriminanzanalyse sind für Datensätze mit vielen Merkmalen weniger geeignet. → Alternative Klassifikationsbäume <ul style="list-style-type: none"> • Diskriminanzanalyse und logistische Regression (insbes. multinomial): Anzahl zu schätzenden Parameter wird schnell zu gross. • Obwohl k-NN keine Modellparameter hat, leidet es unter Fluch der Dimensionalität → im hochdimensionalen Raum werden viel mehr Punkte für gleiche Abdeckung benötigt. Beispiel: Wenn 10 Beobachtungen 10% des Raums in einer Dimension abdecken, sind 100 Beobachtungen in 2 Dimensionen erforderlich. • kNN und Diskriminanzanalyse funktionieren nur begrenzt mit kategorialen Merkmalen. 				

	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
Klassifikationsbäume	<ul style="list-style-type: none"> • Intuitives Verfahren, um Objekte mittels «ja/nein» oder «wahr/falsch» Fragen einer Klasse zuzuordnen. • Nächste Frage hängt von Antwort der Frage zuvor ab. • Fragen/Klassenzuordnung werden im Baum gezeigt • Basis für Ensemble-Verfahren: Bagged Trees, Random Forest, Boosting (Adaboost) • Verwendung auch für Regressionsprobleme • Kann mehrere Bäume geben, die zu Daten passen <p>Einschub: Odds</p> <ul style="list-style-type: none"> • Wahr' als Chancenverhältnisse ausdrücken → Odds • Kennt man aus Umgangssprache oder Wetten • Z.B. Chance stehen 5 zu 3. • Geschrieben als Odds $\frac{5}{3} = 1.7$ • Als Wahr' $\frac{5}{8} = 0.625$, Gegenwahr': $1 - 0.625 = 0.375$ • Odds aus Wahr' $Odds = \frac{p}{1-p} = \frac{0.625}{0.375} = 1.7$ • Gewinnchancen kleiner bei Odds zwischen 0 und 1, Gewinnchancen grösser bei Odds zwischen 1 und ∞ → unsymmetrisch → Alternative log(Odds) <p>logit-Funktion</p> <ul style="list-style-type: none"> • Wahr' zwischen 0 und 1 • Odds zwischen 0 und ∞ • log(Odds) zwischen -∞ und ∞ • Modellierung der log-Odds (logit-Funktion) mit linearen Prädiktor. $\log\left(\frac{p(y x)}{1-p(y x)}\right) = \beta_0 + \beta_1 \cdot x$ • $p \rightarrow p(y x) = \frac{1}{1+\exp(-(\beta_0+\beta_1 \cdot x))} \rightarrow$ Sigmoid Funktion <p>Wo unterteilen?</p> <ul style="list-style-type: none"> • Fragen so stellen, dass Teilmengen möglichst einheitlich sind (d.h. nur Objekte der gleichen Klasse). • Ist zweckmässiger mit Verunreinigung I (impurity) • Welche Verunreinigung an Knoten N in zwei Untergruppen (N_L und N_R) so stark wie möglich verringert • $\Delta I(N) = I(N) - \alpha_L \cdot I(N_L) - (1 - \alpha_L) \cdot I(N_R)$ • mit α_L Anteil Beobachtung in Untergruppe N_L · Häufigste Verunreinigungsmassnahme Gini-Index: Erwartete Fehlerrate Knoten N, falls Label zufällig von den am Knoten vorhandenen Klassen selektiert werden. • $Gini = 1 - \sum_{i=1}^g p_i^2$, p_i: relative Häufigkeit der Klasse i in diesem Knoten · Alternative Entropy $Entropy = \sum_{i=1}^g p_i \cdot \log_2(p_i)$ 	<p>Vorgehen (Baumkonstruktion):</p> <ul style="list-style-type: none"> • Ausgangslage: gesamter Raum • Unterteilungen Raum mit binären Fragen, welche sich auf eine Variable beziehen → nur orthogonale Aufteilungen möglich (→ nur gerade Linien) • Aufteilung bis keine weiteren Unterteilungen mehr möglich sind (auch andere Abbruchkriterien möglich). • Klassenzuordnung gemäss Mehrheit im Endknoten. <p>Stoppregel</p> <ul style="list-style-type: none"> • Baum wachsen bis $\Delta I(N) \leq 0 \rightarrow$ muss nicht reine Knoten geben → häufig überangepasst Bäume (Overfitting) <p>Mögliche Stoppregeln:</p> <ul style="list-style-type: none"> • $\Delta I(N)$ soll Schwellenwert nicht unterschreiten → direkt kontrollierbar in <code>rpart</code> über Parameter <code>cp</code> in <code>caret</code> ist das der Hyperparameter. $Schwellenwert = cp \cdot I(N_{start})$ • Minimale Anzahl Objekte, bei denen ein Knoten noch aufgeteilt wird → <code>train(., method = "rpart", control = rpart.control(minsplit = 20))</code> • Mindestanzahl Objekte, welche in Endknoten noch vorhanden sein müssen → <code>rpart(., method = "rpart", control = rpart.control(minbucket = round(minsplit/3))</code> <p>CP mit KV optimieren! Cost-complexity pruning:</p> <ul style="list-style-type: none"> • Baum zuerst vollständig wachsen lassen • cp mit Kreuzvalidierung optimieren • Dann vollständigen Baum stutzen zu optimalem cp <p>Was ist die richtige Komplexität eines Modells?</p> <ul style="list-style-type: none"> • Bias-Variance Tradeoff: Fehler = Bias + Varianz + un-reduzierbarer Fehler <p>Cost-complexity Pruning (kürzen, stutzen)</p> <ul style="list-style-type: none"> • Abwägung Genauigkeit und Komplexität des Baums: • $R(\alpha) = Fehlerrate + \alpha \cdot T$ • T: Anzahl Endknoten und α: Komplexitätsparameter • Für $\alpha = 0$ ist der grösstmögliche Baum optimal. • Wie finden «optimale» α? Führen K-fache KV für versch. Werte von α für $i = 1, \dots, K$ wie folgt durch: <ol style="list-style-type: none"> 1. Lass Baum vollständig auf Trainingsdaten i wachsen 2. Schneide Knoten Schritt für Schritt ab und berechne $R(\alpha)$ mit den Validierungsdaten i. 3. Speichere optimale $R(\alpha)$. → Mittlere Fehlerquote als Funktion der α • Endgültiger Baum: Gestutzter Baum mit allen Trainingsdaten, der dem optimalen Baum entspricht α. • $cp = \frac{\alpha}{R(\infty)}$, ($R(\infty)$ Cost complexity an Wurzel) 	<ul style="list-style-type: none"> + anschauliche Darstellung, gut für Kommunikation mit Fachleuten und Auftraggebern mit wenig Interesse für statistische Details + Interpretierbarkeit + Keine Verteilungsannahmen + Funktioniert mit kategorialen Daten + Robust gegenüber Ausreissern + Keine Variablentransformation notwendig (oft doch sinnvoll !!!) + Kann nichtlineare Strukturen & lokale Wechselwirkungen erfassen + Erwartungsgetreuer Schätzer + Kleiner Fehler, wenn aussagekräftige Eingabevariablen vorhanden sind und Baum ausreichend tief ist <ul style="list-style-type: none"> - neigt zu Overfitting * hohe Variation → selten für Vorhersagen verwendet - Bäume sind instabil → Die Interpretation kann sich ändern. - Erfordert grosse Datenmengen, um additive Strukturen oder lineare Strukturen zu erfassen. 	<pre>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)) # Visualisierung des Baums plot(r.tree\$finalModel, margin = 0.2, uniform = TRUE) text(r.tree\$finalModel, use.n = TRUE) # Alternative für schöner Visualisierung library(rattle) fancyRpartPlot(r.tree\$finalModel) # Vorhersage für Test-Daten confusionMatrix(predict(r.tree, newdata = testDat), testDat\$Species) # Alternative direkt mit rpart library(rpart) tree <- rpart(Species ~., data = trainDat, control = rpart.control(xval = 10, cp = 0)) printcp(tree) # 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 = 0.01)) plot(tree)</pre>

Bias-Variance Tradeoff			Ensemble-Methoden		
Bias-Variance Tradeoff	Warum Tradeoff?	Umgang Tradeoff	Ensemble-Methoden	Beispiel Ensemble-Methode	Pro (+) / Contra (-)
<ul style="list-style-type: none"> Bias/Varianz wesentlich bei Modellbewertung Bias ist Differenz zwischen mittleren Vorhersage des Modells und richtigen Wert, den wir vorhersagen wollen ($Bias(\hat{Y}) = E(\hat{Y}) - Y$). Modell mit hohem Bias schenkt Trainingsdaten sehr wenig Aufmerksamkeit, was zu Modell führt, welches zu stark vereinfacht ist. Es hat hohen Fehler bei Trainings- und Testdaten. Varianz ist die Variabilität der Modellvorhersage für einen bestimmten Datenpunkt. Modell mit hoher Varianz schenkt Trainingsdaten zu viel Aufmerksamkeit und verallgemeinert nicht auf ungesehene Daten. Wenn Trainingsdaten geringfügig ändern, ändert auch Modell. Modelle mit hoher Varianz schneiden bei Trainingsdaten sehr gut ab, haben aber hohe Fehlerquoten bei ungesehenen Daten. Entscheid.bäume = hohe Varianz, klein Bias 	<ul style="list-style-type: none"> Wenn Modell zu einfach, kann es hohen Bias und geringe Varianz aufweisen. Wenn Modell zu komplex, wird es hohe Varianz und geringen Bias aufweisen. Komplexitätstradeoff ist Grund, warum Tradeoff zwischen Bias/Varianz besteht. Algorithmus kann nicht gleichzeitig komplex und weniger komplex sein. → Müssen richtiges Gleichgewicht zwischen Unter- und Überanpassung finden. $Gesamtfehler = Bias^2 + Varianz + Irreduzierbarer Fehler$ 	<ul style="list-style-type: none"> Modellauswahl: Testen div. Modelle, um optimale Lsg. für Problem zu finden. Durch richtige Wahl finden gutes Gleichgewicht Bias/Varianz. Kreuzvalidierung: Bewertung, wie Modell im Training auf ungesehene Daten reagiert, und Anpassung Hyperparameter. Regularisierung: Beinhaltet Strafterm, um zu verhindern, dass Modell zu sehr an Trainingsdaten anpasst und somit zu Überanpassung führt. Feature engineering: Bearbeitung Eingabefeatures zur Verbesserung Modelleistung. Ensemble-Methoden 	<ul style="list-style-type: none"> Ansatz, um Klassifizierer noch leistungsstärker zu machen Hypothese, dass Kombination mehrerer Modelle oft leistungsfähigeres Modell ergeben kann Viele Instanzen v. gleichen Klassifizierers oder versch. Klassifizierer zusammen verwenden. Gesamtklassifizierer → z.B. Mehrheitsabstimmung aller einzelnen Klassifizierer. Erfolgreich, wenn einzelne Klassifizierer im Schnitt unverzerrt <ul style="list-style-type: none"> Reduktion der Variabilität gegenüber den einzelnen Klassifizierern gut wenn einzelne Klassifizierer stochastisch unabhängig 	<ul style="list-style-type: none"> Bsp. mit 25 Basisklassifizierer Jeder der Klassifizierer hat eine Fehlerrate $\epsilon = 0.35$ Angenommen alle Klassifizierer sind unabhängig (fragwürdigste Annahme) Nehmen Mehrheitsentscheid Mehrheitsentscheid ist falsch, wenn 13 oder mehr der einzelnen Klassifizierer falsch liegen. $X \sim Bin(size = 25, p_0 = 0.35) \Rightarrow P(X \geq 13) = \sum_{i=13}^{25} \epsilon^i (1 - \epsilon)^{25-i} = 1 - pbinom(12, 25, 0.35) = 0.06$ Ensembles sind nur besser als einzelner Klassifizierer, wenn jeder Klassifizierer besser ist als zufälliges Raten! 	<ul style="list-style-type: none"> + Verbesserung Genauigkeit + Macht Modell robuster + Ermöglicht gleichzeitige Erfassen von einfacheren lineare und komplexeren Beziehungen in Daten - Reduziert Modell Interpretierbarkeit - Zeitaufwändiger - Auswahl Modelle für Erstellung Ensembles ist nicht immer einfach und braucht Erfahrung

Ensemble-Ansätze			
Ensemble-Ansätze	Stacking → Beispielcode siehe Seite 14	Idee Bagging	Klassifizierung
<ul style="list-style-type: none"> Stacking: heterogene schwache Klassifizierer werden parallel trainiert und mittels Metamodell kombiniert, um Vorhersage basierend auf Vorhersagen der versch. schwachen Modellen zumache Bagging (bootstrapping & aggregieren): homogene schwache Klassifizierer, lernen parallel anhand unabhängiger Bootstrap-Samples und werden aggregiert. → Bagged Trees und Random Forest (Bagged trees mit «speziellen Trick») Boosting (iteratives Verfahren zur adaptiven Änderung der Trainingsdaten durch stärkere Fokussierung auf zuvor falsch klassifizierte Datensätze). → Gradient boosting, ADABOOST und XGBoost 	<ol style="list-style-type: none"> Daten aufteilen in 3 Datensätze: Train1-, Train2-Daten und Testdaten Training mit verschiedene Algorithmen auf den Train1-Daten (Hyperparameter-Tuning mittels Kreuzvalidierung). → erste Ebene Wenden Modelle auf Train2-Daten an und speichern vorgeschagten Wahrscheinlichkeiten jedes dieser Algorithmen. Trainiere Klassifikator mit vorhergesagten Wahrscheinlichkeiten aus der ersten Schicht (Feature) und ursprünglichen Zielvariable → Toplayer Die Wahrscheinlichkeitsvorhersagen werden für neue Daten (Testdaten) mit den Modellen der ersten Schicht gemacht, die als Eingangsvariablen für das Modell der obersten Schicht verwendet werden. → endgültige Vorhersage (Verschiedene Anpassungen sind möglich: z.B. Mehrheitsentscheidung der ersten Schicht, Neutraining der ersten Schicht nach Training der obersten Schicht mit allen Trainingsdaten, weitere Schichten oder noch komplexere Ansätze) 	<ol style="list-style-type: none"> Mehrere Datensätze anhand originale Trainingsdaten erstellen Erstelle mehrere Klassifizierer Kombiniere Klassifizierer <p>Grundidee (Beispiele Bäume):</p> <ol style="list-style-type: none"> Züchten viele Bäumen auf Bootstrap-Stichproben aus Trainingsdaten Minimiere Fehler, indem man Bäume ausreichend tief wachsen lässt Reduzierung Varianz der Variablen der unverfälschten Bäume durch Mittelwertbildung <p>→ Hochgradig nichtlineare Schätzer wie Bäume profitieren am meisten vom Bagging.</p>	<p>Wie klassifizieren neue Beobachtung?</p> <p>I. Jeder Baum hat «Siegerklasse». Verwende Klasse, die am häufigsten gewonnen hat.</p> <p>II. Mitteln Wahr: $p(c v) = \frac{1}{T} \sum_t p_t(c v)$</p>

Kategoriale Features				
Allgemein	Label Encoding	One Hot Encoding	Target Encoding	Beispiel Target Encoding
<ul style="list-style-type: none"> Nicht alle Modelle können kategoriale Variablen verarbeiten (z.B. k-NN oder DA) Workarounds mit nachfolgenden Methoden. 	<ul style="list-style-type: none"> Jedes Level kategoriale Variable durch Zahl ersetze Bei k-NN ist dies jedoch kritisch (Abstände sind nicht mehr korrekt). Random Forest funktioniert besser, ist aber auch nicht effizient. 	<ul style="list-style-type: none"> besser → kategoriale Variable in mehrere Variablen aufteilen (eine Variable pro Level) + Ermöglicht Modell, die von Haus aus keine kat. Var. verarbeiten können (z. B. k-NN), dies effektiv zu tun; geht von keiner Ordnung zwischen Levels aus - Nachteil: Erhöht Dimensionalität Merkmalsraum, was Modelltraining erschwert; kann Merkmalsauswahl in Random Forests beeinträchtigen. 	<ul style="list-style-type: none"> wenig bekannte, aber wirksame Transformation von kat. Var. Jede Stufe eines kat. Merkmals wird durch durchschn. Antwort in Zielvariablen ersetzt. Vorgehen 1. Gruppieren Trainingsdaten nach Stufen einer kategorischen Variable 2. Berechnen Durchschnitt/Anteil der Zielvariablen für jede Stufe 3. Weisen Durchschnitt jeder Beobachtung zu, die zu dieser Stufe gehört + Geringere Dimensionalität des Merkmalsraums; Erfassung des Einflusses von Kategorien auf Zielvariablen, nützlich in datenreichem Kontext - Ggf. Verzerrung mean durch schlecht repräsentierte oder stark variable Kat. Data leakage: Wenn Zielvariable in Codierung verwenden, kann Overfitting geben. Implementiert in Funktion <code>ranger(..., respect.unordered.factors = "order")</code> 	<p>Feature category Cat und Zielvariable Buy:</p> <pre>td <- data.frame(Cat = c("A", "B", "A", "C", "B", "B", "C", "A"), Buy = c("y", "y", "y", "n", "n", "n", "y", "y"))</pre> <pre>(ratio_buy <- tapply(td\$Buy, td\$Cat, FUN = function(x) mean(x == "y")))</pre> <pre>td\$EncCat <- ratio_buy[td\$Cat]</pre>

One-Hot-Encoding ist ideal für Modelle, die eindeutige, separate kategoriale Merkmale benötigen. Target-Encoding besser für direkte Modellierung Beziehung zwischen Kategorien und Zielvariablen.

	Idee/Ansatz	Voraussetzungen/Vorgehen	Vorteil (+) / Nachteil (-)	R-Code [Default]
Random Forest	<p>Gegeben: N Trainingsdaten und M Features/Merkmale</p> <p>Evaluation Random Forest: Out of Bag (oob) Fehlerrate</p> <ul style="list-style-type: none"> Bei jedem Baum werden nicht verwendeten Testdaten jeweils zur Evaluation der Fehlerrate verwendet. Mitteln über alle Bäume → KV on the fly! <p>Kategoriale Variablen mit Random Forest</p> <ul style="list-style-type: none"> Bäume können i.A. gut mit kategorialen Var. umgehen Bei vielen Faktorstufen wird dies numerisch komplex (Permutation Aufteilungsmöglichkeiten) und langsam. <code>R randomForest</code> kann nur mit 53 Stufen umgehen. One-hot-encoding ist nicht optimale Strategie für RF. Problem: zufällige Variablenauswahl in Knotenpunkten. Effekt: andere Merkmale werden herunter gewichtet und feature importance ist nicht mehr korrekt. In <code>caret</code> wird One-Hot-Codierung automatisch mit Formelnotation durchgeführt! Sogar mit nur ein paar Faktorstufen. (kann dies verhindern, indem ohne Formelschreibweise oder direkt mit <code>randomForest</code> arbeitet. 	<ul style="list-style-type: none"> RF ist nicht immer die beste Methode - aber oft funktioniert sie einfach! <p>Jeder Baum wird mit Algorithmus konstruiert:</p> <ol style="list-style-type: none"> Trainingsset indem N mal mit Zurücklegen aus allen N verfügbaren Trainingsfällen gezogen wird (Bootstrap) → restliche Daten für Validierung Erzeuge ungestutzten Klassifikationsbaum mit Modifikation: Bei jedem Knoten wird beste Aufteilung nur unter zufälligen Auswahl von mtry Features gesucht. mtry $\ll M$ (default Klassifikation mtry = \sqrt{M}) Entkorrelieren mit allen Mitteln! → Maximierung der Varianzreduktion indem Korrelation zwischen Bäumen möglichst klein gehalten wird → Verhindert Overfitting Vorhersage aus Mehrheitsabstimmung aller Bäume 	<ul style="list-style-type: none"> Einfache Handhabung - keine Verteilungsannahmen/Transformationen Perfekte Benchmark-Methode Liefert kostenlose oob Fehlerrate Kann mit tausenden von Eingangsvariablen umgehen, auch wenn viele Störvariablen darunter sind Guter Umgang mit Wechselwirkung Robust gegen Ausreisser Neigt nicht zu Overfitting (Beobachtungen müssen unabhängig sein!) Kann unbalancierte Daten durch Bootstrap-Sampling verarbeiten Anwendung auf Regressionsproblem Kann verwendet werden, um fehlende Daten zu imputieren Bietet Methoden zur Evaluierung der Bedeutung der Variablen und Abhängigkeitsdiagramme Nur 53 Faktorstufen 	<pre>ctrl <- trainControl(method = "oob") r.rf <- train(Species ~., data = trainDat, method = "rf", trControl = ctrl, tuneLength = 3, ntree = 1000) #ntree Anzahl Bäume r.rf\$results plot(r.rf\$finalModel, las = 1) # Konfusionsmatrix r.rf\$finalModel\$confusion # Alternative mit library(ranger) r.out <- ranger(formula = survived ~., data = trainDat, num.trees = 1000, mtry = 2, re- spect.unordered.factors = "order") #"order" = target encoding #sonst level encoding (Standard) confusionMatrix(predict(r.rf, data = testDat), testDat\$survived)</pre>

Unbalancierte Daten

Allgemein	Ansätze Umgang	ROC-Kurve	Anpassung Prior	Sampling Methoden	Random Forest
<p>Haben unbalancierten Datensatz mit mehr Beobachtungen in Klasse 1 als in Klasse 2 (bspw. 85 % F zu 15 % T). Für hohe Genauigkeit, reicht aus, jede Beobachtung Klasse 1 zuzuweisen.</p> <ul style="list-style-type: none"> Sensitivität (True-Positive-Rate) Anteil "Positiven", die korrekt identifiziert werden. Spezifität (True-Negativ-Rate) Anteil "Negativen", die korrekt identifiziert werden. Erhalten in diesem Beispiel hohe Spezifität, aber tiefe Sensitivität. Je nach Fragestellung kann es andersherum sein. 	<p>Was machen, um Identifikation zu verbessern?</p> <p>Alternativer Cutoff</p> <ul style="list-style-type: none"> Wahrscheinlichkeitsvorhersagen und dort alternativen Cutoff-Wert wählen (Standard 50 %), neu z.B. 20 % Cutoff-Wert von Fragestellung/Anforderungen abhängig Auswertung auf Validierungsdaten Trade-of zwischen verschiedenen Arten von Fehlern. → ROC-Kurve <pre>pred_prob <- predict(res.glm, newdata = testDat, type = "prob") pred_neu <- factor(ifelse(pred_prob[,1] > 0.03, "Yes", "No"), levels = c("Yes", "No")) #levels: Klassen, Cutoff neu 3 % confusionMatrix(data = pred_neu, reference = test- Dat\$default)</pre>	<p>ROC-Kurve (Receiver Operator Characteristic): Sensitivity (True-Positive Rate (tpr)) gegen 1-Specificity (False-Positive-Rate (fpr)) für verschiedene Cut-Off-Werte</p> <ul style="list-style-type: none"> Modell mit perfekter Trennung geht durch den Punkt (1, 1) Möglicher Cutoff-Wert nächster Punkt dazu auf der Kurve <p>ROC-Kurve zeigt, dass Modellvergleiche auf Basis Accuracy irreführend sein können, da bessere Cutoff-Wert möglich sind.</p> <p>ROC kann auch zur Modellevaluierung verwendet werden</p> <pre>library(pROC); r_roc <- roc(test- Dat\$Churn, pred_prob[,1], las = 1, levels = c("F", "T")); plot(r_roc, las = 1) (coord <- coords(r_roc, x = "best", best.method = "closest.topleft")) pred_alt <- factor(pred_prob[,1] > coord[,1], levels = c("T", "F")) confusionMatrix(data = pred_alt, refer- ence = testDat\$Churn)</pre>	<ul style="list-style-type: none"> Einige Klassifizierer können Prior-Wahrscheinlichkeit einstellen (Satz Bayes). → Naive Bayes-Klassifikator oder Diskriminanzanalyse Standardmässig wird Prior-Wahrsch' aus Verteilung in Trainingsdaten bestimmen. Kann das auch anpassen. ACHTUNG: Wahrscheinlichkeiten stimmen nicht mehr und dürfen nicht mehr interpretiert werden <pre>lda_fit <- train(Churn ~, method = "lda", prior = c(0.5,0.5)) conf\$table; conf\$by- Class["Sensitivity"]</pre>	<p>Prior Sampling</p> <ul style="list-style-type: none"> Falls Dysbalance bekannt, Datenaufnahme anpassen (balancierte Trainingsdaten, Testdaten sollten für korrekte Evaluation, trotzdem tatsächlichen Verteilung entsprechen) <p>Up-Sampling</p> <ul style="list-style-type: none"> Ziehen aus kleineren Klasse bis Klassen angeglichen sind Gibt ausgeklügeltere Ansätze, z.B. SMOTE (synthetic minority over-sampling technique) <p>Down-Sampling</p> <ul style="list-style-type: none"> Beobachtungen aus Mehrheitsklasse löschen, um Anzahl Minderheitsklasse anzugleichen. (Optimal: Matching Verfahren) → Datenverlust spannend im Kontakt von Bootstrapsamples → Stratifizierte Samplen beim Random Forest <p>Bemerkung: Resampling erst nach Aufteilen Trainings/Testdatensatz durchführen.</p> <pre>down <- downSample(x = trainDat[, -15], y = trainDat\$Churn, yname = "Churn") dim(down); #774 15 table(down\$Churn) # 387 387 rf_d <- train(x = down[, -15], y = down\$Churn, method = "rf", trControl = ctrl) pred_d <- predict(rf_d, newdata = testDat) confusionMatrix(pred_d, testDat\$Churn) up <- upSample(...) # down überall durch up ersetzen</pre>	<p>Down-Sampling mit Random Forest → Stratifiziert Bootstrap-Stichprobe</p> <pre>ctrl <- trainControl(method = "oob") r.rf <- train(x = trainDat[, -15], y = trainDat\$Churn, method = "rf", tuneLength = 3, trControl = ctrl, strata = trainDat\$Churn, sampsize = rep(387, 2)) #strata = Zielvar. #sampsize = Vektor mit länge = Anz. Klassen + Anz. Elemente v kleinster Klasse aus Trainingset table(trainDat\$y)</pre>

Interpretierbares Machine Learning					
Variable Importance VI – Wichtigkeit der Variablen					
Interpretation Modelle	Allgemein	Ansätze	Modellspezifisch Ansätze	Modellunabhängige Ansätze	ICE Plot
<ul style="list-style-type: none"> Machine Learning Algorithmen, insbesondere Ensemble-Verfahren, haben Ruf Black-Box-Modelle zu sein. Das ist irreführend, man kann auch bei komplexen Algorithmen Zusammenhänge der Variablen analysieren (ähnlich wie lineare Regression). Explainable Machine Learning ist aktives Forschungsfeld. Licht in Modelle bringen. <p>Fokus auf 2 Aspekte:</p> <ul style="list-style-type: none"> Wichtigkeit Variablen Einfluss Variablen 	<p>Wie wichtig ist eine Variable bei einem bestimmten Modell, um genaue Vorhersagen zu machen. → je mehr sich Modell auf eine Variable verlässt für Vorhersagen, desto wichtiger ist sie.</p> <p>Wieso ist das wichtig?</p> <ul style="list-style-type: none"> Wollen sehen, ob Modell realistisch ist und Erwartungen entspricht. Haben Möglichkeit Modell zu verbessern? Feature-Selektion (welche Variablen für gutes Modell) 	<p>Modellspezifisch</p> <ul style="list-style-type: none"> Ansatz direkt mit Modell verknüpft Nicht für alle Algorithmen verfügbar Nicht über verschiedene Algorithmen vergleichbar <p>Modellunabhängig</p> <p>Flexible, für alle Klassifikationsalgorithmen verwendbar.</p>	<ul style="list-style-type: none"> <code>caret</code> hat modellspez. Ansätze (Rand. Forest) Bei jedem Split, bei dem Variable verwendet wird, werden Werte der Variablen im Out-of-Bag-Sample zufällig gemischt, wobei alle anderen Variablen gleichbleiben. Gemessen wird Abnahme Vorhersagegenauigkeit. Differenz zwischen den beiden Genauigkeiten wird dann über alle Bäume gemittelt. → Idee auch modellunabhängig umsetzbar <pre>ctrl <- trainControl(method = "oob") r.rf <- train(x = trainDat[, -15], y = trainDat\$Churn, method = "rf", importance = TRUE) varImp(r.rf, scale = FALSE) scale = TRUE skaliert, dass max. Wert bei 100 liegt. train importance = TRUE, wird bei mehr als 2 Klassen Importance für Klassen einzeln ausgegeben, sonst werden diese entsprechend Klassenhäufigkeit gewichtet.</pre>	<ul style="list-style-type: none"> Permutationsansatz v. Random Forest kann auf beliebige Modell verallgemeinert werden (ohne Out-of-Bag-Daten). Trainieren eines beliebigen Algorithmus Loop über alle m Variablen 1. Werte der i-ten Variable in Daten zufällig mischen. Werte der restlichen $m - 1$ Variablen bleiben fix 2. Vorhersage für modifizierter Datensatz. 3. Messen Abnahme Vorhersagegenauigkeit. Ranking basierend auf Vorhersageabnahme bei einzelnen Variablen (je stärker Abnahme, desto wichtiger Variable). <pre>library(vip); q_fit <- train(Churn ~., data = trainDat, method = "qda") r_vi <- vi(q_fit, train = testDat, target = "ZIELVARIABLE", method = "permute", metric = "accuracy", pred_wrapper = predict, nsim = 3) #nsim Anz. Wiederholungen (Zufallsprozess) vip(r_vi, geom = "boxplot")</pre>	<p>Idee: "Flache" Linien haben keinen Einfluss (Variable unwichtig), Linien, die stark ändern (Variable wichtig)</p> <ul style="list-style-type: none"> Veränderungsmass Linie: Varianz Variablen Importance nicht zwingend gleich für verschiedene Ansätze. Variablen Importance bei kategoriellen Variablen nicht immer aussagekräftig. <pre>library(vip) res_vi <- vi(q_fit, train = trainDat, method = "firm", target = "Churn", metric = "accuracy", pred_wrapper = predict, ice = TRUE) #firm für VI vip(res_vi)</pre>
Allgemeine Einschränkung: Alle Ansätze hängen von Modellqualität ab. Wenn Modell schlecht ist, ist Variable Importance nicht aussagekräftig.					
Einfluss einer Variable					
ICE-Plots (Individual Conditional Expectation)			R-Code [Default]		
<ul style="list-style-type: none"> Intuitiver Ansatz, um Einfluss einer Variable auf Klassifizierung zu analysieren (ICE oder Partial Dependence Plot PDP) Jede Beobachtung wird Werte Variable über ganzen Wertebereich variiert, während alle anderen Variablen festgehalten werden, und dafür eine Vorhersage erstellt. → Linie für jede Beobachtung. X-Achse zeigt Wertebereich und y-Achse wie sich Wahrscheinlichkeit für Zielklasse ändert. Rote Linie zeigt Mittelwert → Partial Dependence Plot (PDP) ICE-Plot ist reinen PDP vorzuziehen, da sich dort gegenläufig Effekt ausmitteln können und verschwinden. Vorsicht Wechselwirkungen zwischen Variablen werden ICE-Plot nicht berücksichtigt. Kann falsche Schlussfolgerungen 			<pre>library(pdp); r_par <- partial(q_fit, pred.var = "ERKL. VARIABLE", type = "classification", which.class = "WERT ZIELKLASSE", prob = T, ice = T) plotPartial(r_par, main = "ICE Plot") pred.var: Welche erklärende Variable soll genommen werden? type: Auch für Regressionsprobleme verwendbar which.class: Welche Klasse der Zielvariable soll gezeigt werden? prob: Wahrscheinlichkeit, sonst logit-Werte, ice: ICE, sonst PDP-Kurve</pre>		

Eigenes Ensemble-Modell mit Stacking erstellen

```
intrain <- createDataPartition(y = churn$Churn, p = 0.8, list = F)
train <- churn[intrain,]; test <- churn[-intrain,]
intrain2 <- createDataPartition(y = train$Churn, p = 0.7, list = F)
train1 <- train[intrain2,]; train2 <- train[-intrain2,]

ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3)
m_knn <- train(Churn ~., data = train1, method = "knn", trControl = ctrl, tuneLength = 10)
m_lda <- train(Churn ~., data = train1, method = "lda", trControl = ctrl)
m_qda <- train(Churn ~., data = train1, method = "qda", trControl = ctrl)
m_log <- train(Churn ~., data = train1, method = "glm", family = "binomial", trControl = ctrl)
results <- resamples(list(KNN = m_knn, LDA = m_lda, QDA = m_qda, LOG = m_log))
summary(results); bwplot(results); dotplot(results) #Modellvergleich

pred_knn <- predict(m_knn, train2, type = "prob")[, 1];
pred_lrg <- predict(m_log, train2, type = "prob")[, 1]
pred_lda <- predict(m_lda, train2, type = "prob")[, 1]
```

```
pred_qda <- predict(m_qda, train2, type = "prob")[, 1]
input_toplayer <- data.frame(pred_knn = pred_knn, pred_lrg = pred_lrg,
                             pred_lda = pred_lda, pred_qda = pred_qda)

#Verwenden Wahr' aus ersten Layers als Feature für Top Layer mit einer log Regression
input_toplayer$Churn <- train2$Churn
(model_top <- train(Churn ~., data = input_toplayer, method = 'glm', trControl = ctrl))

pred_knn <- predict(m_knn, test, type = "prob")[, 1] #Mit Testdaten auswerten
pred_lrg <- predict(m_log, test, type = "prob")[, 1]
pred_lda <- predict(m_lda, test, type = "prob")[, 1]
pred_qda <- predict(m_qda, test, type = "prob")[, 1]
input_toplayer <- data.frame(pred_knn = pred_knn, pred_lrg = pred_lrg,
                             pred_lda = pred_lda, pred_qda = pred_qda)
confusionMatrix(predict(model_top, newdata = input_toplayer), test$Churn)
```

Boosting					
Allgemein	AdaBoost	Gradient Boosting	Verlustfunktion	Stochastic Gradient Boosting	XGBoost Extreme Gradient Boosting
<ul style="list-style-type: none"> Im Unterschied zum Bagging, werden verschiedene homogenen Klassifizierer nicht mehr unabhängig voneinander, sondern sequenziell trainiert. Grösseres Gewicht auf Beobachtungen, welche im vorherigen Modell falsch vorhergesagt wurden. Unterschiedliche Gewichtung Modelle im Ensemble. Hauptaugenmerk auf Reduzierung Bias → Basismodelle, Modelle mit geringer Varianz, aber hohem Bias (Unterschied zu Bagging!). Diese Modelle sind weniger rechenaufwendig. Modelle meist einfache, flache Bäume 	<ul style="list-style-type: none"> Erste praktische Implement. Binärer Klassifikator 3 Hauptideen Kombination von vielen schwachen Lerner → Bäume; häufig mit nur einer Stufe (stamps). Nicht alle Modelle sind gleich gut. Bessere Modelle erhalten im Ensemble mehr Gewicht. Jedes Modell berücksichtigt die Fehler vom vorgängigen stärker. <p>Berücksichtigung Gewichte in Klassifizieren</p> <ul style="list-style-type: none"> Gewichteter Gini-Index (d.h. stärkere Berücksichtigung der Beobachtungen mit grösserem Gewicht) Neues Datensample → ziehen mit Zurücklegen aus Originaldaten mit Gewicht entsprechen den Wahr' → Neuer Klassifiziere mit Sample (gleiche Gewichtung) <pre>ctrl <- trainControl(method = "cv", number = 5) res.ada <- train(label ~ ., data = trainDat, method = "AdaBoost.M1", trControl = ctrl, tuneGrid = data.frame(mfinal = 50, maxdepth = 1, coeflearn = "Breiman")) confusionMatrix(predict(res.ada, newdata = testDat), reference = testDat\$label) mfinal: Anzahl Bäume maxdepth: max Tiefe einzel Bäume coeflearn: Versch. Ansätze Gewichtung Ensembles (α)</pre>	<ul style="list-style-type: none"> Schätzen Zielvariable mit schwachen Algorithmus (Regression/Klassifikation) Schrittweises pushen Vorhersage gegen beobachtete Zielvariable anhand Informationen aus Residuen. $\hat{y} = F_0(x) + \sum_{m=1}^M \lambda \cdot f_m(x) = F(x)$ \hat{y}: Vorhersage $F_0(x)$: Initialisierungsmodell für Zielvariable $\sum_{m=1}^M \dots$: Kumulative Verfeinerung Modell mit M Iterationen $f_m(x)$: Weak learner (bei Klassifikation einfache Klassifikationsbäume) mit Residuen (gradient) aus Schritt $m - 1$ λ: Lernrate zur Steuerung Schrittweite $F(x)$: finales Modell <p>Algorithmus:</p> <ul style="list-style-type: none"> Initialisierung $F_0(x)$ ist i.d.R. konstanter Schätzer. Mittelwert bei Regression. Frequenzen bei Klassifikation. Funktionen $f_m(x)$ verwenden i.d.R. kleine Regressionsbäume. 1. Initialisieren Vorhersagemodell: $\hat{y} = F_0(x)$ 2. Von $m = 0$ bis M <ol style="list-style-type: none"> Residuen berechnen: $r = y - F_m(x)$ Residuen modellieren: $\hat{r} = f_m(x)$ Neues Modell: $F_{m+1}(x) = F_m(x) + \lambda f_m(x)$ <p>+ Sehr leistungsfähig, flexibel - Viele Hyperparameter - Anfällig auf Overfitting - Interpretierbarkeit</p>	<ul style="list-style-type: none"> Verlustfunktion definiert, was Fehler zwischen vorhergesagten Werten \hat{y} und tatsächlichen Werten y angesehen wird. Quantifiziert, wie gut Modell $F(x)$ der Zielvariable annähert. Vorhersagen völlig falsch → Verlustfunktion hohen Wert, wenn gut → niedrigere. Variabel anwendbar; z.B. Regression (kontinuierliche Zielvariable) oder Klassifikation (kategoriale Zielvariable). Ziel → Verlustfunktion minimieren, z.B. mit Algorithmus. Regression: <ul style="list-style-type: none"> $MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$, $MAE = \frac{\sum_{i=1}^n y_i - \hat{y}_i }{n}$ Klassifikation: <ul style="list-style-type: none"> 0-1 Verlust $Genaugigkeit = \sum_{i=1}^n \begin{cases} 1: y_i \neq \hat{y}_i \\ 0: y_i = \hat{y}_i \end{cases}$ Negativ Log Likelihood (NLL)/Cross-Entropie für Wahr'vorhersage \hat{p} von M Klassen: $NLL = \sum_{i=1}^n (-\sum_{c=1}^M y_{i,c} \log(\hat{p}_{i,c}))$ <p>Verbesserung Vorhersage</p> <ul style="list-style-type: none"> Schauen für jeden Datenpunkt, in welche Richtung verbessern könnten. Vorhersagen werden in diese Richtung "gedrückt". Praxis: wissen nicht, wo Optimum liegt. Nutzen Gradienten, um Richtung zu bestimmen. Für Vorhersage steht Gradienten nicht zur Verfügung. $f(x)$ modelliert Gradienten in Abhängigkeit Daten. $-Gradienten \approx f(x)$ können für zu schätzenden Gradienten Residuen $r_i = y_i - \hat{y}_i$ einsetzen! Klassifikation Log Loss $\mathcal{L}(\hat{y}_i) = y_i \log(\hat{p}_i) \rightarrow \frac{d\mathcal{L}}{d \log(\hat{o} \hat{a} s_i)} = -y_i + \hat{p}_i = r_i$ 	<ul style="list-style-type: none"> Stabilere Optimierung und geringere Gefahr «Feststecken» lok. Optima. Algorithmus: 1. Initial. Vorhersagemodell: $\hat{y} = F_0(x)$ 2. Von $m = 0$ bis M <ol style="list-style-type: none"> Residuen berechnen: $r = y - F_m(x)$ und zufälliges Ziehen daraus (Bspw. 50 %) Residuen modellieren: $\hat{r} = f_m(x)$ Neues Modell: $F_{m+1}(x) = F_m(x) + \lambda f_m(x)$ <p>Parameter:</p> <ul style="list-style-type: none"> M: Anz. Bäume (1000-10000) λ: Lernrate (0-1); kleine Rate reduziert Effekt einzelner Bäume, verbessert aber Vorhersage auf lange Sicht (braucht mehr Bäume). Kleinere Werte liefern bessere Performance. Mit kleineren Werten steigen Computerressourcen. Empfehlung 0.01-0.001. Baumtiefe: i.A. kleine Bäume, Empfehlung 2-5. Min Anz. an Beobachtungen in Endknoten der Bäume (z.B. 10) Gefahr Gradient Boosting: Neigt schnell zu Overfitting, da kein Strafterm! <pre>ctrl <- trainControl(classProbs = TRUE, method = "cv", number = 5, summaryFunction = mnLogLoss) gbm_fit <- train(Spe ~., data = trainDat, method = "gbm", trControl = ctrl, distribution = "multinomial", metric = "logLoss", verbose = FALSE, tuneGrid = data.frame(expand.grid(n.trees = c(200, 500, 1000), interaction.depth = seq(2, 5), shrinkage = 0.01, n.misnobsnode = 10))) #Verteilung "bernoulli"/"gaussian" confusionMatrix(predict(gbm_fit, newdata = testDat), testDat\$Spe)</pre>	<ul style="list-style-type: none"> erfolgreiche Implementation von Gradient Boosting + Bietet Regularisierungsoptionen, um Modellkomplexität zu reduzieren → verringert Overfitting. + Berücksichtigt bei Baumentwicklung reduzierte Anzahl Features (vergleichbar zum Random Forest). + Betracht Gradienten 2. Ordnung → hilft besser Richtung zur Minimierung Lossfunktion zu finden. <p>Parameter:</p> <ul style="list-style-type: none"> M (nrounds): Anz. Bäume, mit CV optimieren. Baumtiefe (max_depth): Wie komplex einzelne Bäume max werden. 2-8 λ (eta): Learning rate, wie gross sind die Schritte (wie gbm), 0.1-0.001, 0.01 CP (gamma): Regularisierung Baumkomplexität, (vgl. Entscheidbaum), 0-0.1 Variablenanteil (colsample_bytree): Wie viel Prozent der Variablen pro Baum genutzt werden min Samples je Node (min_child_weight): Mindestanzahl Beobachtungen in einer Partition Subsample: Prozent zufälligen Beobachtungen, die für Erstellung Baum benutzt werden. <pre>ctrl <- trainControl(classProbs = TRUE, method = "none", summaryFunction = mnLogLoss) #braucht one-hot encoding xgb_fit <- train(Species ~., data = trainDat, 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)) confusionMatrix(predict(xgb_fit, newdata = testDat), testDat\$Species)</pre>

Modellevaluierung 2.0			
Allgemein	AUC (Area Under ROC Curve)	Güte einer Wahrscheinlichkeit (Wahr' = Wahrscheinlichkeit)	
		Kalibrierungskurve	Log Loss
<p>Accuracy</p> <ul style="list-style-type: none"> + Leicht zu verstehen und zu interpretieren; direkt anwendbar. - Irreführend bei unbalancierten Klassen – berücksichtigt keine Wahrlichkeiten • Accuracy ist nicht immer gutes Mass für Performanceevaluation Modell. • Gibt bessere Masse vorhergesagten Wahr' nutzen: AUC und Log-Loss • Wahrscheinlichkeit anstelle von Klassifikation: Beim Umwandeln von Wahr' zu Klassen verlieren wir Information! 	<ul style="list-style-type: none"> • ROC-Kurve: (Receiver Operator Characteristic)-Kurve: Sensitivity (tpr) gegen 1-Specificity (fpr) für verschiedene Cut-Off-Werte • AUC (Area under the Curve) Modellbeurteilung anhand der Fläche unter der Kurve (AUC) → zeigt, wie gut Modell Klassen unterscheiden kann (unabhängig Schwellenwert). Höherer AUC-Wert → besseres Modell + Effektiv bei unbalancierten Klassen + berücksichtigt Rangfolge Vorhersagen - Nicht empfindlich gegenüber den absoluten Werten der Vorhersage - Irreführend sein bei sehr unbalancierten Datensätzen 	<p>Statt Wahr' mittels Grenzen in Klassen umzuwandeln, kann Güte der geschätzten Wahr'keiten direkt beurteilt werden.</p> <ul style="list-style-type: none"> • Grafik, die beurteilt, ob Modell gut kalibriert ist. Modell ist gut kalibriert, wenn vorhergesagte Wahrscheinlichkeit mit tatsächlichen Häufigkeit Ereignisse, die es vorhersagt, übereinstimmen. D.h. wenn Modellvorhersage 80 % Wahr' zutrifft, sollte Ereignis, das es vorhersagt, tatsächlich in ca. 80 % der Fälle eintreten. • Kalibrierungskurve vergleicht geschätzte Wahrscheinlichkeiten mit echten Frequenz in beobachteten Daten. • Wahrscheinlichkeiten werden in Gruppen eingeteilt und mit den jeweiligen Mittelwerten gearbeitet. • Vorhergesagte Wahrscheinlichkeiten der Klasse m, $\hat{p}(y = Klasse\ m x)$, werden in Gruppen eingeteilt (Standard = 10 Gruppen) • x-Achse sind gemittelten, vorhergesagten Wahrscheinlichkeiten der Beobachtungen pro Gruppe. • y-Achse sind Anteile Beobachtungen in Gruppen, die zur Klasse m gehören. • Perfekte Kalibrierung: Alle Punkte liegen auf Winkelhalbierenden (oberhalb = Modell überschätzt W'keiten, unterhalb = Modell unterschätzt W'keiten). 	<p>Metrik, welche Güte der Kalibrierung misst. (häufig verwendeten Metriken für Klassifikation) Auch als Cross-Entropie oder negative LogLikelihood bekannt.</p> <ul style="list-style-type: none"> • Gibt an, wie nahe Vorhersagewahrscheinlichkeit an tatsächlichen Klasse liegt. Je mehr Vorhersagewahrscheinlichkeit abweicht, desto höher ist Log-Loss-Wert. • Log-Loss ist datenspezifisch! Zwei Modelle, die auf verschiedenen Datensätze angewendet werden, können nicht verglichen werden. • N = Anzahl Beobachtungen, M = Anzahl Klassen • y_{ij}: 1 wenn Beobachtung zur Klasse j gehört, sonst 0 • \hat{p}_{ij}: vorhergesagte Wahrscheinlichkeit, dass Beobachtung zur Klasse j gehört • $LogLoss = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(\hat{p}_{ij})$ + Sinnvoll für Wahr'schätzungen und nützlich für Optimierung - Schwer zu interpretieren und ausreiseranfällig
<p>ROV Kurve Allgemein</p> <p>Werte liegen auf Winkelhalbierenden = Modell bringt keine Erkenntnis Wenn Werte unterhalb liegen, dann wurde Zuordnung (T/F) vertauscht</p>	<pre>ctrl <- trainControl(classProbs = TRUE, method = "cv", number = 10, summaryFunction = multiClassSummary) iris.knn <- train(Species ~., data = trainDat, method = "knn", trControl = ctrl, tuneLength = 3, metric = "AUC")</pre> <p>ROC-Kurve Vergleich zweier Modelle</p> <pre>library(pROC) #levels/Klassen im Dat sind hoch und tief res.glm <- train(crime ~., data = trainDat, method = "glm", family = "binomial") ctrl <- trainControl(method = "none") r.tree <- train(crime ~., data = trainDat, method = "rpart", trControl = ctrl, tuneGrid = data.frame(cp = 0)) glm_pred <- predict(res.glm, testDat, type = 'prob') tree_pred <- predict(r.tree, testDat, type = 'prob')</pre> <pre>glm_roc <- roc(testDat\$crime ~ glm_pred\$'hoch') tree_roc <- roc(testDat\$crime ~ tree_pred\$'hoch') glm_roc\$auc; tree_roc\$auc #für Erhalt AUC Werte</pre> <pre>plot(glm_roc, las = 1) lines(y = tree_roc\$sensitivities, x = tree_roc\$specificities, col = "red", lwd = 2, lty = 2) legend("bottomright", legend = c(paste("glm auc:", round(glm_roc\$auc, 2)), paste("tree auc:", round(tree_roc\$auc, 2))), lwd = 2, lty = c(1, 2), col = c("black", "red"), bty = "n")</pre>	<pre>library(predtools) calibration_plot(data, obs = 'y', pred = 'prob')</pre>	<pre>ctrl <- trainControl(classProbs = TRUE, method = "cv", number = 10, summaryFunction = multiClassSummary) iris.knn <- train(Species ~., data = trainDat, method="knn", trControl = ctrl, tuneLength = 3, metric = "logLoss")</pre> <p>Log Loss Vergleich zweier Modelle</p> <pre>library(pROC) #levels Dat hoch und tief res.glm <- train(crime ~., data = trainDat, method = "glm", family = "binomial") ctrl <- trainControl(method = "none") r.tree <- train(crime ~., data = trainDat, method = "rpart", trControl = ctrl, tuneGrid = data.frame(cp = 0)) glm_pred <- predict(res.glm, testDat, type = 'prob') tree_pred <- predict(r.tree, testDat, type = 'prob')</pre> <pre>glm_pred\$obs <- testDat\$crime tree_pred\$obs <- testDat\$crime</pre> <pre>mnLogLoss(glm_pred, lev = levels(testDat\$crime)) mnLogLoss(tree_pred, lev = levels(testDat\$crime))</pre>
<p>Welche Metrik sollten Sie verwenden? ⇒ Verwenden Sie mehrere Metriken, um ein umfassendes Bild der Modellperformance zu erhalten.</p>			

Ablauf Klassifizierung

- Datenaufbereitung: Kontrolle Balancierung Daten, Zielvariable als Faktor hinterlegt mit der richtigen Reihenfolge (Faktorstufe von Interesse an erster Stelle)?
- `Crime$crime <- factor(Crime$Crime, levels = c("tief", "hoch"))`
- Variablenkategorien? Numerisch, Factor, binär??? → korrekt hinterlegt
- Verteilung numerische Variablen? → Skalieren??? `boxplot(dat)`
- `pairs(trainDat[, 2:4], col = c("red", "blue")[trainDat$default])` #Korrelation zwischen den einzelnen Variablen als Plot darstellen

Feature Engineering				
Allgemein	Textdaten	Bilddaten	Zeitreihen	Mehrere Messungen
<ul style="list-style-type: none"> Prozess, in welchem Rohdaten in Features umgewandelt werden, die in Algorithmen des maschinellen Lernens, z. B. in Vorhersagemodellen, verwendet werden können. Je nach Format Rohdaten gibt es unterschiedliche Techniken. Um gut Feature erstellen zu können, braucht man Fachwissen und muss die Daten verstehen! Feature-Engineering hat i.d.R. massiv größeren Einfluss auf Qualität Ergebnisse als Algorithmus, welchen man verwendet! ML Projekten Datenaufbereitung und Feature-Engineering rund 80% der Zeit beanspruchen. 	<ul style="list-style-type: none"> Tokenisierung → Text aufteilen Stoppwörter entfernen → Uninformative Wörter entfernen Stemming oder Lemmatisierung → Wörter basierend auf Regeln oder auf Struktur Sprache (Wörterbücher) auf Stammform zurückführen Document-Term Matrix → Tabelle mit Häufigkeiten der Wörter Entfernung weiterer seltene oder sehr häufige Wörter. Bag-of-Words → Vektor mit Anzahl Beobachtungen für jedes Wort in Wortliste <p>Erweiterungen:</p> <ul style="list-style-type: none"> Transformation des Bag-of-Words in TF-IDF word2vec: Ansatz, um automat. Zusammenhänge zwischen Wörtern und Bedeutung aus Textkorpus zu lernen. Token (Wörter/Unterwörter/Multiwort-Ausdrücke) lassen sich als Zahlen ausdrücken, um sie für Modelle zu verwenden. Vortrainierte Transformer-Modelle 	<ul style="list-style-type: none"> Für Klassifikation Bildern wird Deep Learning DL verwenden, welche fähig sind relevante Merkmale automatisch zu lernen. Vortrainierte Modelle nutzen, um Feature aus Bildern zu extrahieren. Feature sind Vektoren, welche für klassische ML Modelle verwenden. 	<ul style="list-style-type: none"> Frequenzen extrahieren mittels Fourier-Transformation. Können DL Methoden, welche Feature selbstständig lernen sehr erfolgreich sein. 	<ul style="list-style-type: none"> Hat man zu einer Beobachtung mehrere Messungen der gleichen Variable (zum Beispiel kurze Zeitreihe) kann es sinnvoll sein, Kennzahlen (min, max, Anzahl Peaks, Median, mean) zu extrahieren und diese als Feature für die Modellierung zu verwenden. I.d.R. je mehr Domainwissen für Feature Engineering verwenden, umso besser.

Fehlende Daten					
Allgemein	Arten von fehlenden Daten	Umgang mit fehlenden Daten	Imputation fehlende Daten	Nützliche Modelle für Imputation	Multiple Imputation
<ul style="list-style-type: none"> sind häufiges Problem bei Data Science Projekten. Gründe sind vielfältig und häufig nicht unter unserer Kontrolle fehlende Daten müssen bei Datenaufbereitung adressieren, da Modelle damit i.d.R. nicht umgehen können. Gibt keine richtige Handhabung, aber Möglichkeiten, die je Anwendungsfall verwenden können. → Wenn fehlende Werte mit falschen Daten auffüllen, fügen wir einen Bias hinzu. 	<ul style="list-style-type: none"> Missing completely at Random (MCAR): keinen Zusammenhang zwischen Fehlen Datenpunkts mit restlichen Werten im Datensatz. Fehlende Daten sind zufällige Teilmenge. Missing at Random (MAR): Fehlende Werte haben Zusammenhang mit beobachteten Daten, aber nicht mit fehlenden Wert selbst (z.B. Umfrage über psychische Störungen bei Männern und Frauen. Männern melden Depressionsstatus vielleicht weniger, das hat aber nichts mit dem Grad ihrer Depressionen zu tun). Missing Not at Random (MNAR): Beziehung zwischen fehlendem Wert und dessen Wert (Z.B. Depressionsumfrage; Befragte mit höheren Depressionswerten füllen Umfrage aufgrund ihres Depressionsgrades häufiger nicht aus). 	<p>Ansatz: Complete Case Analysis (nur vollständige Beobachtungen verwenden)</p> <ul style="list-style-type: none"> Entfernen von Beobachtungen (Zeilen): <code>train(..., na.action = na.omit)</code> Wenn einzelnen Features (Spalten) sehr viele Daten fehlen (z.B. 60%) kann es auch Sinn machen ganze Spalte zu entfernen. In den Fällen MCAR und MAR ist dies in der Regel okay. Im Fall MNAR kann dies zu Bias im Modell führen. Sollte versuchen zu ergründen, wieso Werte fehlen. I.A. ist Weglassen Daten nicht beste Lösung, da wir potenziell nützliche Informationen verlieren. ⇒ Besser zu imputieren. 	<p>Imputation: Statistische Verfahren, um fehlende Daten in Datenmatrix zu vervollständigen. Aktuelles Forschungsfeld. Möglichkeiten:</p> <ul style="list-style-type: none"> Einsetzen konstanten Wertes, der sich von allen anderen Werten unterscheidet (z.B. neue Faktorstufe für NAs bei kategorialen Variable) Einsetzen Mittelwert, Median oder Moduswert für diese Spalte (nicht zu empfehlen) <ul style="list-style-type: none"> Mittelwert: nicht optimal, zerstört Struktur, reduziert Varianz Zufallswert: zerstört Struktur, Beibehalt Varianz Nächstem Nachbar: Benachbarte Punkte müssen über andere, vorhandene Variablen gefunden werden. Schätzen Wert mit eigenen prädiktiven Modell Mehrfache Imputation 	<ul style="list-style-type: none"> KNN: kNN im Paket VIM. Imputation basierend auf Variation Gower-Distanz für numerische, kategoriale, geordnete und semikontinuierliche Variablen. Random Forest: z.B. missForest im Paket missForest: geeignet um gemischte Daten mit komplexen Interaktionen und nichtlinearer Beziehungen zu imputieren. Funktion preProcess() für Imputation: <ul style="list-style-type: none"> <code>"knnImpute"</code> kNN-Klassifizierer → Daten werden automatisch zentriert und skaliert <code>"bagImpute"</code> Bagged trees <code>"medianImpute"</code> Imputation mit Median Trainingsdaten Ansatz muss vorgängig ans Training durchgeführt werden: <pre>imp <- preProcess(data, method = "knnImpute") data_miss_imp <- predict(imp, data)</pre> 	<ul style="list-style-type: none"> Warnung: Imputation werden so behandelt, als seien es Messungen → Unterschätzung der Unsicherheit ⇒ Multiple Imputation (mehrere Datensätze, welche zum Trainieren Modell verwendet werden, analysieren der Varianz). Idee Erzeugen mehrerer imputierter Datensätze Passen Modell an jeden Datensatz an Zusammenführen der Ergebnisse der verschiedenen Modelle Beispiel: MICE (Multivariate Imputation by Chained Equations): Imputiert jede Variable stochastisch unter Verwendung eines Regressionsmodells. R-Paket mice

	Idee/Ansatz	Voraussetzungen/Vorgehen	Pro (+) / Contra (-)	R-Code [Default]
Support Vector Machine (SVM)	<p>• Flexibles Klassifikationsverfahren mit guter Performance. Bspw. geeignet, wenn Anzahl Features sehr gross oder Anzahl Beobachtungen limitiert ist.</p> <p>Grundidee:</p> <ul style="list-style-type: none"> • Jede Beobachtung → p-dimensionaler Vektor • SVM konstruiert Hyperebene (Hyperplane), um 2 Klassen zu trennen. <p>Maximale Margin Klassifikator</p> <ul style="list-style-type: none"> • Maximale Margin Klassifikator wählt Ebene, welche Daten so trennt, dass Abstand von Hyperebene zum nächsten Trainingspunkt maximal ist. • kleine Margin → Gefahr, dass Trainingsdaten überanpasst werden (Overfitting) • grosse Margin → Erwartung, dass auch Testdaten korrekt abgebildet werden <p>Support Vektors</p> <ul style="list-style-type: none"> • Support Vektor: Trainingsdaten, welche nächsten an Trenngrenze liegen, bestimmen Hyperebene • Restliche Trainingsdaten tragen nicht zur Festlegung der Grenze bei! • Haben einen Abstand von $\leq M$ zur Hyperebene <p>Hyperplane bestimmen</p> <ul style="list-style-type: none"> • Hyperplane definiert: $\beta_0 + \beta_1 \cdot X_1 + \dots + \beta_p \cdot X_p = 0$ • Punkte $0 \rightarrow \beta_0 + \beta_1 \cdot O_1 + \dots + \beta_p \cdot O_p > 0$ und $O^* \rightarrow \beta_0 + \beta_1 \cdot O_1^* + \dots + \beta_p \cdot O_p^* < 0$ liegen auf unterschiedlichen Seiten der Hyperebene ⇒ Notationstrick: Zielvariable mit $y = \pm 1$ definieren. <p>Optimierungsproblem:</p> <ul style="list-style-type: none"> • $y_1 \cdot (\beta_0 + \beta_1 \cdot x_1 + \dots + \beta_p \cdot x_p) \geq M$ für alle $i = 1, \dots, N$ • Gesucht β-Vektor, welcher M maximiert. Für eindeutige Lösung benötigen wir Nebenbedingung: $\sum_{j=1}^p \beta_j^2 = 1$ da ansonsten jede Multiplikation β's mit Konstante ein Lsg <p>Nicht eindeutig trennbare Daten</p> <ul style="list-style-type: none"> • Vielfach sind Daten mit linearen Grenze nicht trennbar. • Idee: Erlauben einige Fehlklassifizierungen (Support Vector Classifier) $y_i \cdot (\beta_0 + \beta_1 \cdot x_{i1} + \dots + \beta_p \cdot x_{ip}) \geq M \cdot (1 - \epsilon_i)$, ϵ sind Slack-Variablen, welche Beobachtungen erlauben auf falschen Seite Margin oder Hyperebene zu liegen. • $\epsilon_i \geq 0$, $\sum_{i=1}^n \epsilon_i \leq \frac{1}{C} \Rightarrow C > 0$ (Cost): Tuningparameter, entspricht Kosten Fehlklassifikation in Trainingsdaten 	<p>Voraussetzungen/Vorgehen</p> <p>Kleines C:</p> <ul style="list-style-type: none"> • Mehr Fehlklassifikationen da geringe Kosten → ermöglicht grössere Margin • Gefahr Übergeneralisierung Trainingsdaten • Wenn Trainingsdaten leicht ändern, neigen Grenzen dazu, stabil zu bleiben • Klassifikator mit geringer Variabilität und grossem Bias <p>Grösseres C:</p> <ul style="list-style-type: none"> • Weniger Fehlklassifikationen da grössere Kosten → kleinere Margin • Gefahr Überanpassung an Trainingsdaten • Wenn die Trainingsdaten leicht ändern, können Grenzen stark variieren • Klassifikator mit grosser Variabilität und kleinem Bias ⇒ mehr Fehlklassifikationen zulassen, dafür aber grössere Margin zu erhalten. <p>Nichtlineare Grenzen</p> <ul style="list-style-type: none"> • Manchmal lineare Hyperebene nicht funktionieren, egal was C gewählt wird. <p>Feature Erweiterung Idee:</p> <ul style="list-style-type: none"> • Feature Raum vergrössern, d.h. weitere erklärende Variablen hinzufügen, indem Transformationen der ursprünglichen Variablen als zusätzliche Variablen verwendet werden (z.B. $X_1^2, X_1^3, X_1 X_2, X_1 X_2^2$) → Übergang von p- zum M-dimensionalem Raum ($M > p$). • Anwendung SVM-Klassifizierer auf vergrösserten Raum → ergibt nicht-lineare Entscheidungsgrenzen in originalen Dimension. <p>Kernel Trick ⇒ SVM</p> <ul style="list-style-type: none"> • Polynome manuell hinzufügen, wird schnell schwierig handelbar. → Bessere Lösung, um Nichtlinearität einzubinden sind Kernels. • Kernel → Funktion, die Ähnlichkeit zwischen 2 Beobachtungen quantifiziert. • Ansatz Skalarprodukt $\langle x_i, x_{i'} \rangle = \sum_{j=1}^p x_{ij} \cdot x_{i'j}$ • Alternative Schreibweise für Linear SV Klassifizierer $f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$ Um Parameter α_i und β_0 zu schätzen braucht man $\binom{n}{2}$ Skalarprodukte zwischen allen Paaren an Trainingsbeobachtungen, wobei meisten $\hat{\alpha}_i$ 0 sind. Relevant ist Supportset S (Support Vectors). • Daraus SV Klassifizierer: $f(x) = \beta_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle$ <p>SVM: Radialer Kernel</p> <ul style="list-style-type: none"> • $K(x_i, x_{i'}) = \exp\left(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2\right)$ • γ: Hyperpara, definiert Einflussbereich Trainingsdaten • kleines γ grosser Einflussbereich (Grenzen linearer) • grosses γ kleiner Einflussbereich (flexible Grenzen) • → Schätzung durch Kreuzvalidierung <p>SVM: Mehr als 2 Klassen:</p> <ul style="list-style-type: none"> • OVA: One versus All: Fit K verschiedene SVM-Modelle: jeweils Klasse k gegen den Rest. Klassifiziere eine neue Beobachtung in die Klasse mit der grössten Distanz zur Hyperebene. • OVO: One versus One: Fit alle $\binom{K}{2}$ paarweise SVM-Modelle. Klassifiziere Beobachtung in Klasse, welche meisten paarweisen Vergleiche gewinnt. <p>→ Welche Ansatz wählen? Falls K nicht zu gross → OVO</p>	<p>+ Klassifikationsmethode mit hoher Vorhersageperformance (Erweiterung auf Regressionsprobleme möglich).</p> <p>+ Funktioniert auch mit relativ wenigen Trainingsdaten und vielen Features.</p> <p>+ Auch bei mehr als 2 Klassen anwendbar</p> <p>+ Kernel-Trick (neu konstruierte Features) erlaubt es, nichtlineare Trenngrenzen zu finden.</p> <p>- Separierung Klassen durch Hyperebene und Margin basierend auf SV. Übrige Trainingsdaten tragen nichts zur Hyperebene bei.</p> <p>- Liefert standardmässig keine Wahrscheinlichkeitsvorhersagen.</p>	<p>Linearer Kernel</p> <pre>ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3) svm_lin <- train(y ~ ., data = trainDat, method = "svmLinear2", trControl = ctrl, tuneGrid = data.frame(cost = c(0.1, 1))) svm_lin\$bestTune #Ausgabe C</pre> <p>Bemerkung: Daten werden in Funktion defaultmässig standardisiert. Was auch wichtig ist.</p> <pre>confusionMatrix(predict(svm_lin, testDat), reference = testDat\$y)</pre> <p>Gaussian bzw. radialer Kernel</p> <pre>ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3) svm_rad <- train(y ~ ., data = trainDat, 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</pre> <pre>confusionMatrix(svm_rad, "none")</pre> <p>Polynomieller Kernel</p> <pre>svm_poly <- train(y ~ ., data = trainDat, method = "svmPoly", tuneGrid = expand.grid(degree = 1:3, scale = c(0.1, 1, 2, 5, 10), C = c(0.1, 1, 5, 8, 10, 15)), trControl = ctrl) svm_poly\$bestTune plot(svm_poly)</pre> <pre>confusionMatrix(data = predict(svm_poly, testDat), reference = testDat\$y)</pre>

Deep-Learning DL				
Was neu an DL?	Einordnung Modelle	Power of Deep Learning	FCNN	
<ul style="list-style-type: none"> Traditionelle ML extrahieren handgefertigter Merkmale und verwenden Merkmale zum Modell Training. Deep-Learning DL (End-to-End-Ansatz) Deep-Neuronale Netze starten mit Rohdaten und lernen während Training, geeignete hierarchische Merkmale zu extrahieren, um diese für Klassifizierung zu verwenden. 	<ul style="list-style-type: none"> Regelbasierte Systeme: Input → handdesigntes Programm → Output Klassisches Maschine Learning: Input → handdesigntes Features → Mapping von Features → Output. Extrahieren von handgefertigten Features, welche zum Modelltraining (z.B. SVM, RF) verwendet werden. Representation/Feature Learning (z.B. Deep Learning): Input → Features → Mapping von Features → Output Deep Learning: Ausgehend Rohdaten lernt Modell selbstständig, geeignete Feature zu extrahieren, welche der Vorhersage dienen. Für grössere Aufgaben ist oft grössere Rechenkapazität notwendig (GPU/Server-Lösung). 	<ul style="list-style-type: none"> Ziel neuronalen Netzes ist Funktionen $y = f(x)$ finden, die Attribute x auf Ausgabe y abbilden. Wieso funktionieren das? → Universal Approximation Theorem: Sei f beliebige kontinuierliche Funktion auf $[0,1]^d$ und φ eine beliebige sigmoidale Aktivierungsfunktion, dann lässt sich f uniformverteilt durch endliche Summe Form $\hat{f}(x) = \sum_{j=1}^m w_j^{(2)} \varphi(w_j^{(1)T} x + w_{j,0}^{(1)})$ approximieren. Neuronales Netzwerk mit mindestens einem Hidden Layer und einer ausreichenden Anzahl von Neuronen ist in der Lage jede stetige Funktion auf einem kompakten Intervall mit einer beliebigen Genauigkeit anzunähern. Theorem macht keine Aussage über Effizienz Training oder Anzahl der benötigten Neuronen. Es besagt lediglich, dass ein solches Netzwerk theoretisch existiert. In der Praxis kann es jedoch Herausforderungen geben, die mit Overfitting oder anderen Faktoren zusammenhängen können. 	<p>Fully Connected Neural Networks (FCNN)</p> <ul style="list-style-type: none"> Neuronale Netze: Modell aufbauend auf einer Netzwerkstruktur mit künstlichen Neuronen. Simple Neural Network besteh aus einem Hidden Layer, während Deep Learning Neural Network mehrere Hidden Layers hat. Grosser Aufschwung, da fähig, grosse Mengen unstrukturierter Daten (z.B. Bilder oder Text) besonders gut auswerten zu können. 	
Klassifikationsnetzwerk	Netz trainieren	Ablauf Training	Lernrate	Overfitting
<ul style="list-style-type: none"> w: Gewicht, b: Bias x: Inputvektor $z_1 = b + \sum_{i=1}^n x_i \cdot w_{i1} = b + x_1 \cdot w_{11} + x_2 \cdot w_{21} + \dots + x_n \cdot w_{n1}$ Output \hat{y} letzten Layers wird durch Softmax-Funktion in Wahrsch'vektor \hat{p} für verschiedenen Klassen transformiert: $\hat{p}_i = \frac{e^{\hat{y}_i}}{\sum_{j=1}^m e^{\hat{y}_j}}$ (Summe 1) m = Anzahl Klassen Nicht-lineare Aktivierungsfunktion: <ul style="list-style-type: none"> Sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ (s-förmig) ReLU $R(z) = \max(0, z)$ (lineare Gerade) 	<ul style="list-style-type: none"> Gewichte Netzwerk sind Modellparameter → Deep Learning Modelle haben viele Parameter (Moderne Netzwerke können Billionen (10^9) Gewichte haben.) Anpassung Gewichte durch Optimierung Loss Funktion mit Gradientenverfahren. Steigung Lossfunktionen sind lokale Gradienten. (Ableitung) Iterativer Prozess Schrittweises trainieren mit zufälligen Sample aus Trainingsdaten (Batch). Wird wiederholt, bis alle Trainingsdaten einmal durch sind (Epoche). Repetiert für mehrere Epochen → Je mehr Beispiele Netz für Training bekommt und je öfter es diese gesehen hat, desto besser wird die Performance. 	<ul style="list-style-type: none"> Zufällige Indizierung Gewichte für erste Runde. Forward Pass: Input-Werte werden ins Netz eingespielen und durch Netz propagiert → Model-output <ul style="list-style-type: none"> Vergleich Output mit bekannten Ergebnissen des Datensatzes. Verlustfunktion (z.B. Klassifikation: Log-Loss/Kreuzentropie) $Loss = -\frac{1}{N} \sum_{i=1}^N y_i (\log(\hat{y}_i))$ Backward Pass: Jedes Gewicht und jeder Bias Term wird kleines Stückchen in Richtung angepasst, die Fehler kleiner macht. (→ Gradientenverfahren) Grösse Anpassung wird über Lernrate (Hyperparameter) gesteuert. Ist Rate zu gross, kann Minimum Verlustfunktion verpasst werden. Ist Rate zu klein, dauert Training lange). Forward und Backward Pass werden als iterativer Prozess wiederholt. 	<ul style="list-style-type: none"> ist ein wichtiger Parameter im Deep Learning. hat enormen Einfluss auf Modellanpassung. Deep Learning Modell kann optimal auf Problem zugeschnitten sein und dennoch kann alles schief gehen. → zu gross oder zu kleine Lernrate 	<ul style="list-style-type: none"> Wenn Netz während Training alle Trainingsdaten sehr oft gesehen hat, kann sein, dass es Daten eher auswendig lernt, statt abstraktes Konzept zu lernen. → Overfitting Um abstrahierendes Modells zu erstellen, werden Training Daten in Training und Validation aufgeteilt. Die Gewichte werden jeweils nur auf Trainingsdaten angepasst. Optimierung folgt auf Trainings- und Validierungsdaten. → Finaler Test an ungesehenen Testdaten, welche gar nicht für Training verwendet wurden. Dropout: Während Training werden z.B. 50 % Neuronen eines Layers einfach abgeschaltet. Die Abschaltung erfolgt pro Durchlauf jeweils auf zufälliger Basis. Alle Neuronen des Layers werden damit gezwungen, weniger spezielle und mehr abstrakte Konzepte zu lernen, um trotz reduzierten Anzahl gut zu performen. Für finale Modell (für Vorhersagen) stehen schlussendlich wieder alle Neuronen zur Verfügung. Ergebnisse werden i.d.R. robuster und besser im Hinblick auf neue unbekannte Daten.

Keras Workflow					
Workflow	1. Definieren Netzwerk	2. Kompilieren Netzwerk	Preprocessing Daten	3. Fitten Netzwerk	4. Evaluation Netzwerk
1. Definieren Netzwerk → Layerweise Definition Netzwerk	<code>library(keras3)</code> <code>model <- keras_model_sequential()</code> <code>layer_dense(model, units = 16, input_shape = c(4), activation = 'relu')</code> <code>layer_dropout(model, rate = 0.2)</code>	<code>model <- compile(model, optimizer = optimizer_adam(learning_rate = 0.001), loss = 'categorical_crossentropy', metrics=c('accuracy'))</code> #AUC	<code>library(caret)</code> <code>index <- createDataPartition(iris\$Species, p = 0.7, list = F)</code> <code>index <- sample(index)</code> <code>trainDat <- as.matrix(iris[index, 1:4])</code> <code>testDat <- as.matrix(iris[-index, 1:4])</code> <code>preP <- preProcess(trainDat, method = c("center", "scale"))</code> <code>train_scale <- predict(preP, trainDat)</code> <code>test_scale <- predict(preP, testDat)</code> <code>train_y <- as.numeric(iris[index, 5])-1</code> <code>test_y <- as.numeric(iris[-index, 5])-1</code> <code>test_labels <- to_categorical(test_y)</code> <code>head(train_labels, n = 2)</code>	<code>history <- fit(model, x = train_scale, y = train_labels, batch_size = 32, validation_split=0.2, verbose=0)</code> <code>plot(history)</code>	<code>evaluate(model, x = testDat, y = test_labels, verbose = 0)</code>
2. Kompilieren Netzwerk → Loss und Optimierer definieren	<code>layer_dense(model, units = 16, activation = 'relu')</code> <code>layer_dropout(model, rate = 0.2)</code> <code>layer_dense(model, units = 3, activation = 'softmax')</code> <code>summary(model)</code>	• optimizer: Adaptive Moment Estimation (Adam; <code>optimizer_adam</code>) oder Stochastic Gradient Descent (SGD; <code>optimizer_sgd</code>). • learning_rate: Lernrate. • loss: <code>categorical_crossentropy</code> entspricht Log-Loss. Bei Regression <code>mean_squared_error</code>	<code>train_scale <- predict(preP, trainDat)</code> <code>test_scale <- predict(preP, testDat)</code> <code>train_y <- as.numeric(iris[index, 5])-1</code> <code>test_y <- as.numeric(iris[-index, 5])-1</code> <code>test_labels <- to_categorical(test_y)</code> <code>head(train_labels, n = 2)</code>	• epochs: Epochen zum Trainieren Modell • batch_size: Anzahl Stichproben pro Gradienten-Update • validation_split: Anteil Daten, welche nur Validieren und nicht zum Fitten verwendet werden. • verbose: 0 Silent, 1 Progressbar, 2 eine Linie pro Epoche.	5. Anwendung <code>predict(model, testDat[1:10,])</code>
3. Fitten Netzwerks → Trainieren Modell mit Trainingsdaten	• units: Anzahl Neuronen im Layer • input_shape: Dimension Inputs (nur im ersten Layer notwendig) • activation: Aktivierungsfunktion (relu, sigmoid oder softmax) • rate: Anteil Auszulassende Gewichte • Alternativ Modellaufbau <code>%>% dplyr</code> .		• <code>preP <- preProcess(trainDat, method = c("center", "scale"))</code> • <code>train_scale <- predict(preP, trainDat)</code> • <code>test_scale <- predict(preP, testDat)</code> • <code>train_y <- as.numeric(iris[index, 5])-1</code> • <code>test_y <- as.numeric(iris[-index, 5])-1</code> • <code>test_labels <- to_categorical(test_y)</code> • <code>head(train_labels, n = 2)</code>		
4. Evaluieren Netzwerk mit Testdaten			• Input benötigt numerische Matrizen • Empfehlung Daten skalieren (verhindert Stabilitätsproblem, sorgt effiziente Gewichtanpassung, Regularisierungseffekt) • One Hot Encoding für kat. Daten notwendig (Zielvariable).		
5. Anwendung finales Modell					

Netzwerk Architekturen	Beispiel CNN-Layer	Convolution Layer mit mehreren Filtern	Bilden CNNs in keras
<ul style="list-style-type: none"> • Bis jetzt auf Fully Connected Neuronale Netzwerke (FCNN) beschränkt. • Diese eignen sich für tabellarische Daten. Reihenfolge Eingabe spielt keine Rolle. • Geordnete Daten (z.B. Bilder, Zeitreihen) werden «geflattet». Geordnete Information geht verloren. → Gibt bessere Alternativen. • Für Bilder sind Convolutional Neural Networks (CNNs) geeignet, (verg. ImageNet-Competition) • Während FCNN Gewichte für jedes einzelne Pixel haben, werden bei CNNs Filter mit geteilten Gewichten verwendet, welche über das ganze Bild gleiten. • Weitere Architekturen: z.B. RNNs (Rekurrentes neuronales Netz) bzw. LSTMs (Long short-term memory) für Zeitreihen. 	<ul style="list-style-type: none"> • An jeder Position des Inputs werden die gleichen Gewichte verwendet. • Soll die Größe des Bildes erhalten bleiben, kann z.B. Null-Padding angewendet werden (Rand wird um 0-Werte erweitert). • Maxpooling: Maximum in Filtergröße; reduziert die Größe der Bilder → weniger Gewichte. 	<p>Kernell mit einem Kanal</p> <p>Filter mit mehr als einem Inputkanal</p> <p>Typisches klassisches CNN-Modell</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> <ul style="list-style-type: none"> • filters: Anzahl Filters • kernel_size und pool_size: Filtergröße

Tricks bei der Anwendung
<ul style="list-style-type: none"> • Deep Learning braucht nicht viel Preprocessing, aber Standardisierung Daten hilft oft. Achtung die Aktivierungsfunktion muss zum Wertebereich der Daten passen. • Haben nie genügend Daten für Modelltraining → Data Augmentation. Dabei werden, um Datenmenge zu vergrößern, leicht modifizierte Kopien von bereits vorhandenen Daten erstellt (Bild rotieren). • Was kann man bei limitierten Daten, Computerressourcen und Zeit machen? → Verwenden von vortrainierten Modellen zur Merkmalsgenerierung. <ul style="list-style-type: none"> ○ Neu trainiert wird nur der/die letzte/n Layer/s (Transfer Learning). Alternative mit einfacherem parallelen Modell (Side Learning). ○ Verwenden eines klassischen Klassifikator (SVM) zur Klassifikation basierend auf den extrahierten Features. ○ Insbesondere für Bilderkennung stehen viele vortrainierte Modelle zur Verfügung. Verwendet werden sollte ein Modell, das für einen möglichst ähnlichen Task entwickelt wurde.

Stand der Arbeit

SW	Vorlesung	AB/Aufgabe	R-Skript	ZF
1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
4	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
5	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
6	<input type="checkbox"/> Ostern	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
7	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
8	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
9	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
10	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
11	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
12	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
13	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
14	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Modulaufgaben

Obligatorische Leistungsnachweise für die Modulbewertung:

- Praktische Arbeiten, 40%
- Endprüfung (90 Minuten), 60%
- Die Modulendnote setzt sich somit zusammen aus: $0.4 \times$ Praktische Arbeiten + $0.6 \times$ Endprüfung.
- Für die praktische Arbeit müssen 4-mal im Semester die erlernten Methoden auf einen eigenen Datensatz angewendet und in einem kurzen Bericht dokumentiert und interpretiert werden. Die Berichte werden bewertet (Gewichtung jeweils 8%). Zudem muss einmal im Semester eine vorgestellte Methode inklusive Übungsaufgabe zusammengefasst und der Klasse präsentiert werden (Gewichtung 8%). Die Tasks können in 2er-Gruppen bearbeitet werden. Details zu den Semestertasks sind auf Moodle und in den Vorlesungsunterlagen ausgeführt.

SEP-Hilfsmittel

- Erlaubte Hilfsmittel für Endprüfung: Eine zehenseitige Zusammenfassung auf Papier und ein (elektronisches) File mit R / Python Befehlen. Laptop mit R und/oder Python. Alle benötigten Bibliotheken müssen installiert sein. Es sind keine weiteren Hilfsmittel und elektronischen Kommunikationsmittel wie z.B. Mobiltelefone und WLAN erlaubt.

Statistisches Data Mining – Zusammenfassung

Allgemeines

- Data Mining ist die Anwendung geeigneter Methoden zur (automatisierten) **Entdeckung** von **Strukturen** und **Zusammenhängen** in **grossen Datenmengen**.
- Die Datenbestände werden dabei nach Regelmässigkeiten, Mustern und Strukturen, Abweichungen, jeglicher Art von Zusammenhängen und gegenseitigen Beeinflussungen untersucht.
- Data Mining ist durch den Begriff Data Science in den Hintergrund gerückt. Die Tätigkeit Data Mining ist allerdings nur ein Teil von Data Science und wird dort oft unter dem Begriff Data Analytics subsummiert.
- Data Mining hat grosse Ähnlichkeiten zur statistischen Datenanalyse und verwendet u.a. auch Methoden des Maschinellen Lernens.

Data Mining

Ziele:

- Daten verstehen/visualisieren, Muster finden und klassifizieren → Querverbindungen/Trends erkennen.
- Vorhersagen machen → Fokus liegt auf guten Vorhersagen. Das Modell ist im Allgemeinen nicht so wichtig.
- Hypothesen formulieren.

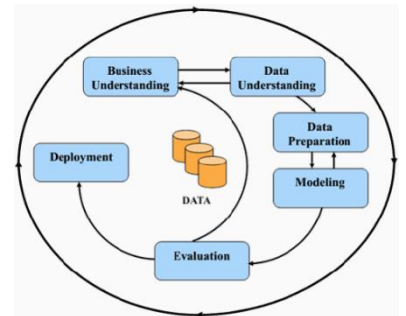
Spezielle Merkmale:

- Anwendung statistischer Methoden auf **grosse Datenbestände** («Big Data»).
- Von einer Beobachtungseinheit (z.B. ein Patient) können zig-tausend Merkmale erhoben werden, so dass die **Anzahl der Merkmale** die **Anzahl der Beobachtungseinheiten übersteigen** kann.
- Charakteristisch, dass Daten **nicht** für bestimmten **Zweck** - z.B. Überprüfung Hypothese - erhoben wurden.
- **Methoden anwendbar, wenn traditionelle Techniken versagen** (d.h. Datenmenge zu gross, zu heterogen oder Dimensionalität zu hoch)

Data Mining Prozess

Data Mining Prozess für datengetriebene Lernen. Cross-industry standard process data mining (CRISP-DM) ↓

- Business Understanding: Was will der Auftraggeber? Wer sind die Empfänger? Was sind die Ressourcen?
- Data Understanding: Verfügbarkeit, Meta-Daten, Struktur
- Data Preparation: Umgang mit fehlenden Werten und Ausreissern, Standardisierung, Transformation, Aggregieren und Selektion von Variablen.
- Modeling: Auswahl des Data Mining Verfahrens, Prognosen, Auswertungen
- Evaluation: Statistische Evaluierung des Modells, qualitative Beurteilung
- Deployment: Umsetzung in die Praxis
- Modellierung/Evaluation stehen im Kurs im Zentrum, auch wenn Data Understanding/Preparation 80 % Zeit ausmachen.



Daten sind der Schlüssel

Prozess

- Business Understanding
- Data Understanding
- Data Preparation
- Modeling
- Evaluation
- Deployment
- Data Mining / Machine Learning lernt direkt von Daten. Um es gut zu machen, brauchen i.d.R. viele Daten!

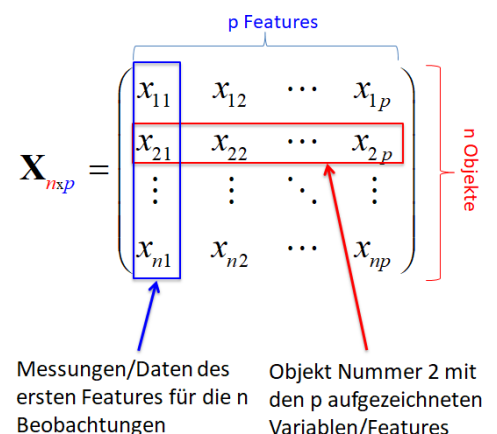
Was kann falsch gehen?

→ garbage in, garbage out (GIGO)

- Problem nicht verstanden
- Falsche Daten
- Fehlende Daten
- Umgang mit grossen Datenmengen
- Falsche Methoden / Annahmen
- Kommunikation

Allgemeine Struktur der Daten

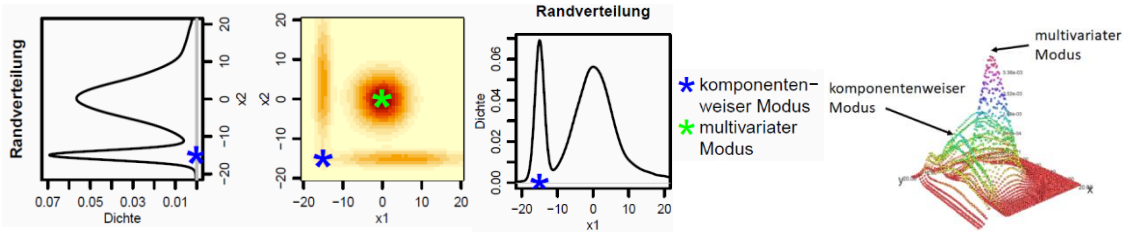
- Datenobjekte und deren Attribute (Flat Table wie Excel Sheet)
- Variable/Features sind Eigenschaften eines Objekts
 - Beispiele: Geschmackskategorie (Whisky), Transaktionsmerkmale, Wörter, Pixel, etc.
 - Für Klassifizierung gibt es ein ausgezeichnetes Attribut (Klassenattribut): z.B. Betrug ja oder nein
- Eine Sammlung von Attributen beschreibt ein Objekt
 - Objekt auch: Individuum, Fall, Datum oder Rekord
 - Beispiele: Whisky, Transaktion, Mail, Bild
- Teilweise müssen die Daten zuerst in die entsprechende Struktur gebracht werden → Feature Engineering



Fallstricke bei univariaten Analysen

- Wieso wollen wir multivariable Daten analysieren? Was ist das Problem, wenn man die Variablen eines multivariaten Datensatz einzeln betrachtet?

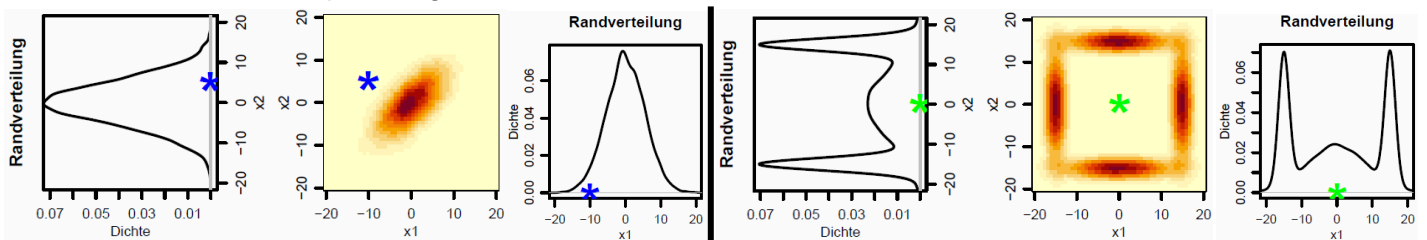
- Beispiel Modus: Multivariater Modus entspricht wahrscheinlichsten Realisierung p-dimensionalen Merkmalsvektors.



- Der multivariate Modus kann nicht über die Moduswerte der Randverteilung abgeleitet werden!
- Beim Median gibt es nicht einmal eine allgemein akzeptierte Definition für den multivariaten Median.
- Die weitest verbreitete Version für den multivariaten Median ist der **L1-Median**. Er ist definiert als Punkt mit minimaler Summe der absoluten Abstände zu allen anderen Punkten (**L1median**, Paket **robustX**).
- Der komponentenweise Median ergibt nicht dasselbe.
- Immerhin ist Mittelwert gutmütig. Der komponentenweise Mittelwert entspricht multivariaten Schwerpunkt.

Ausreisser im multivariaten Raum

- Ein multivariater Ausreisser ist nicht unbedingt auch ein Ausreisser in den Randverteilungen!
- Ausreiser in blau respektive grün in den Grafiken.



Curse of dimensionality (Fluch der Dimensionalität)

- Wie viele Daten werden für eine bestimmte Abdeckung benötigt?
- Bei uniform-/gleichmässig verteilten Daten gilt zudem:
 - In hochdimensionalen Räumen «extreme Merkmalswerte» nicht selten
 - Kaum vorstellbar, aber in hohen Dimensionen wahr: Alle Punkte haben fast gleichen Abstand zueinander

Dimension	Anz. Beobachtung
1	10
2	100
3	1000
...	...
6	10^6
p	10^p