

Formeln				
Lineare Regression	Logistische Regression	Neuronale Netzwerke	Entscheidungsbäume	Reinforcement Learning
<p>Linearer Zusammenhang zwischen den Eingabevariablen x und der Ausgabevariable y wird modelliert.</p> <p>Hypothesenfunktion: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$</p> <p>Kostenfunktion (MSE): $J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$ mit n... Anzahl Trainingsdaten Je kleiner $J(\theta)$, desto besser die Hypothese.</p> <p>Ziel: Finde Parameter θ um J zu minimieren min $J(\theta)$</p> <p>Multivariat: Mehrere Features x_1, x_2, \dots, x_n</p> <p>Polynom-Regression: $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots$ Wird eingesetzt, wenn wenn Features nicht-linearen Zusammenhang haben.</p>	<p>Hypothese: $h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$ $g(z) = \frac{1}{1+e^{-z}}$..... Sigmoidfunktion mit $\theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$</p> <p>Klassifikation: $h_{\theta}(x)$ ist quasi Wahrscheinlichkeit (0...1) $h_{\theta}(x) \geq 0.5 \rightarrow$ Klasse 1 $h_{\theta}(x) < 0.5 \rightarrow$ Klasse 0</p> <p>Entscheidungsgrenze: $\theta^T x = 0$ weil da wird $h_{\theta}(x) = 0,5$</p> <p>Nicht-linearität: $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \dots)$</p> <p>Polynom-Regression: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_{2+} \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots$ Fähigkeit, nicht-lineare Entscheidungsgrenzen zu erzeugen</p>	<p>$z_j^{(l)}$ ist der gewichtete Input für Neuron j in Layer l. $a_j^{(l)}$ ist die Aktivierung des Neurons j in Layer l.</p> <p>Parametermatrizen: $\dim(\theta^{(l)}) = (s_{l+1} \times (s_l + 1))$ θ^l.....Gewichtsmatrizen, l..... layer, s_l..... Anzahl Units in Layer (l), in $(s_l + 1)$, das +1 steht für Bias-Unit (zusätzliche Spalte)</p> <p>Fehlerterm δ: $\delta_{ij}^{(l)} = h_{\theta}(x^{(i)}) - y^{(i)}$ i... Trainingsbeispiel, j... Neuron, l... layer, $h_{\theta}(x^{(i)})$... Netzwerkvorhersage, $y^{(i)}$... tatsächlicher Zielwert, n... Anzahl Trainingsbeispiele,</p> <p>Gradientenberechnung (Backpropagation): $\frac{\partial}{\partial \Theta_{ij}^{(l-1)}} J(\Theta) = \frac{1}{n} \sum_{k=1}^n \delta_{ki}^{(l)} a_j^{(l-1)}$</p> <p>Feedforward: $z^{l+1} = \theta^l a^l$ $a^{l+1} = g(z^{l+1})$</p> <p>Backpropagation: $\delta_j^{(l)} = \left(\sum_{k=1}^{s_{l+1}} \left(\Theta_{kj}^{(l)} \cdot \delta_k^{(l+1)} \right) \right) \cdot g' \left(z_j^{(l)} \right)$ $\Theta_{kj}^{(l)}$... Gewicht zwischen Neuron j in Layer l und Neuron k im Layer $l + 1$ $\delta_k^{(l+1)}$... Fehlerterm aus der folgenden Schicht (bereits berechnet). $g' \left(z_j^{(l)} \right)$ Ableitung der Aktivierungsfunktion am Wert $z_j^{(l)}$ $\delta^L = a^L - y$ $\delta^l = (\theta^l)^T \delta^{l+1} \cdot g'(z^l)$</p> <p>Gradientenabstieg: $\theta^l := \theta^l - \alpha \delta^l a^{l-1}$</p> <p>Aktivierungsfunktionen: Sigmoid, Tanh, ReLU, Leaky ReLU, Softmax</p>	<p>Grundidee: Baumstruktur zur schrittweisen Entscheidung basierend auf Features.</p> <p>Split-Kriterien:</p> <ul style="list-style-type: none"> Entropie: $H(S) = - \sum p_i \log_2(p_i)$ Informationsgewinn Gini-Index: $\text{Gini}(S) = 1 - \sum p_i^2$ <p>Vorteile:</p> <ul style="list-style-type: none"> Interpretierbar Keine Skalierung notwendig <p>Nachteile:</p> <ul style="list-style-type: none"> Overfitting bei tiefen Bäumen Instabil bei kleinen Datenänderungen <p>Pruning (Beschneiden): Reduziert Komplexität und Overfitting</p> <p>Ensemble-Methoden:</p> <ul style="list-style-type: none"> Random Forests: viele Bäume, Voting Boosting (z.B. AdaBoost): sequentielle Optimierung 	<p>Grundidee: Ein Agent lernt durch Interaktion mit einer Umgebung, um Belohnungen zu maximieren.</p> <p>Zentrale Begriffe:</p> <ul style="list-style-type: none"> Agent, Environment Zustand s, Aktion a, Belohnung r, Politik $\pi(a s)$ Ziel: Maximiere erwarteten kumulierten Reward <p>Belohnungsformel: $R_t = \sum_{k=0}^{\infty} \gamma^k r_{\{t+k+1\}}$ mit Diskontfaktor $\gamma \in [0, 1]$</p> <p>Q-Learning: $Q(s, a) := Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$</p> <p>Strategien:</p> <ul style="list-style-type: none"> Exploration vs. Exploitation var ε-greedy Policy <p>Anwendungen:</p> <ul style="list-style-type: none"> Spiele (z.B. AlphaGo) Robotik Empfehlungssysteme
Gradient Descent	Support Vector Machines	Principial Component Analysis(PCA)		Modell Evaluation
<p>Update-Regel: $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ α..... Lernrate $\frac{\partial}{\partial \theta_j} J(\theta)$..... Ableitungsterm, gibt Richtung des steilsten Anstiegs an</p> <p>Update-Regel für lineare Regression: $\theta_j := \theta_j + \alpha \frac{1}{n} \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x_j^{(i)}$ mit θ_j als Parameter und $x_0^{(i)} = 1$ für den Bias-Term.</p> <p>Lernrate α: Zu groß \rightarrow Divergenz, zu klein \rightarrow (zu) langsame Konvergenz</p> <p>Eigenschaften:</p> <ul style="list-style-type: none"> Bei geeigneter Lernrate konvergiert Gradient Descent zu einem lokalen Minimum der Kostenfunktion. Bei konvexen Funktionen sogar zum globalen Minimum. Eine konstante Lernrate kann ausreichen, wenn sie angemessen gewählt ist. Adaptive Lernraten können die Konvergenz verbessern. Steilere Kostenfunktionen erfordern kleinere Lernraten 	<p>Ziel: $\min_{w,b} \left\{ \frac{1}{2} \ w\ ^2 + C \sum_{i=1}^n \xi_i \right\}$</p> <p>Nebenbedingungen: $y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i$, mit $\xi_i \geq 0 \quad (i = 1, \dots, n)$</p> <p>C kontrolliert Trade-off: großes $C \rightarrow$ weniger Fehler, kleines $C \rightarrow$ größerer Margin</p> <p>Kernel: $K(x, l) = \exp \left(-\frac{\ x-l\ ^2}{2\sigma^2} \right)$</p> <p>Kernel-Funktionen:</p> <ul style="list-style-type: none"> Linear: $K(x_i, x_j) = x_i^T x_j$ Für lineare Trennungen Polynomial: $K(x_i, x_j) = (x_i^T x_j + r)^d$ r..... konstanter Offset d..... Grad des Polynoms Erlaubt gekrümmte trennungen (Parabeln, ...) RBF (Gaussian): $K(x_i, x_j) = e^{-\gamma \ x_i - x_j\ ^2}$ Komplexe, nicht-lineare Trennungen Sigmoid: $K(x_i, x_j) = \tanh(\kappa x_i^T x_j + c)$ Inspiriert von neuronalen Netzen, selten verwendet 	<p>Ziel: Reduktion der Dimensionalität bei maximalem Erhalt der Varianz.</p> <p>Schritte:</p> <ol style="list-style-type: none"> Zentrieren der Daten Kovarianzmatrix berechnen Eigenvektoren & -werte berechnen Hauptkomponenten auswählen (größte Eigenwerte) Projektion der Daten auf neue Achsen <p>Mathematisch: Gegeben Datenmatrix X, berechne $C = \frac{1}{n} X^T X$ Finde Eigenvektoren v mit $Cv = \lambda v$</p> <p>Eigenschaften:</p> <ul style="list-style-type: none"> Unüberwachtes Verfahren Hauptachsen sind orthogonal Keine Label nötig <p>Anwendungen:</p> <ul style="list-style-type: none"> Visualisierung Vorverarbeitung für ML Rauschreduktion 		<p>Konfusionsmatrix: TP = True Positive - Patienten, die krank sind und als krank klassifiziert wurden FP = False Positive - Patienten, die gesund sind, aber als krank klassifiziert wurden TN = True Negative - Patienten, die gesund sind und als gesund klassifiziert wurden FN = False Negative - Patienten, die krank sind, aber als gesund klassifiziert wurden</p> <p>Metriken: Accuracy = $\frac{TP + TN}{TP + TN + FP + FN}$ Precision (Relevanz) = $\frac{TP}{TP + FP}$ Recall (Sensitivity) = $\frac{TP}{TP + FN}$ F1-Score = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$</p> <p>Beispiel: Ein Modell mit hohem Recall identifiziert möglichst viele positive Fälle, auch auf Kosten von mehr False Positives.</p> <p>Wann ist Recall wichtiger?</p> <ul style="list-style-type: none"> Krankheitserkennung Sicherheitschecks Betrugserkennung

Aufgaben				
Regularisierung	Backpropagation Aufgabe (21)			
<p>Kostenfunktion mit L2-Regularisierung: $J(\theta) = \frac{1}{2n} \left(\sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^d \theta_j^2 \right)$</p> <p>Effekt von λ (Regularisierungsparameter):</p> <ul style="list-style-type: none"> $\lambda = 0 \rightarrow$ kein Penalty kleines $\lambda \rightarrow$ schwache Bestrafung von großen Parameterwerten, Gefahr von Overfitting großes $\lambda \rightarrow$ starke Bestrafung von großen Parameterwerten, Gefahr von Underfitting <p>Bias-Term θ_0 wird oft nicht regularisiert</p>	<p>Gegeben:</p> <ul style="list-style-type: none"> $x^{(1)} = (x_1^{(1)}, x_2^{(1)}) = (0.5 \ 0.7),$ $y^{(1)} = (1),$ (“Klasse 1”), $h_{\theta}(x^{(1)}) = 0.5,$ Aktivierungsfunktion: Sigmoid $a_j^{(i)} = g\left(z_j^{(i)}\right) = \frac{1}{1+e^{-z_j^{(i)}}}$ Gewichtsmatrizen: $\Theta^{(1)} = \begin{pmatrix} 1 & 1 & 2 \\ 2 & 1.5 & 0.7 \\ -0.5 & -1 & 2 \end{pmatrix},$ $\Theta^{(2)} = \begin{pmatrix} 1 & -0.9 & -0.7 & 0.9 \end{pmatrix}$ <p>Forward-Propagation:</p> <ul style="list-style-type: none"> Schritt 1 - Input Layer (inklusive Bias): $a^{(1)} = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} \end{pmatrix} = \begin{pmatrix} 1 & 0.5 & 0.7 \end{pmatrix}$ Schritt 2 - Hidden Layer Input berechnen: $z^{(2)} = \Theta^{(1)} \cdot a^{(1)} = \begin{pmatrix} 1 & 1 & 2 \\ 2 & 1.5 & 0.7 \\ -0.5 & -1 & 2 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0.5 & 0.7 \end{pmatrix}$ $= \begin{pmatrix} 1 \cdot 1 + 1 \cdot 0.5 + 2 \cdot 0.7 & 2 \cdot 1 + 1.5 \cdot 0.5 + 0.7 \cdot 0.7 & -0.5 \cdot 1 - 1 \cdot 0.5 + 2 \cdot 0.7 \end{pmatrix}$ $= \begin{pmatrix} 2.9 & 3.24 & 0.4 \end{pmatrix}$ Schritt 3 - Aktivierungsfunktion (Sigmoid) anwenden: $a^{(2)} = g(z^{(2)}) = \left(1, \frac{1}{1+e^{-z_1^{(2)}}}, \frac{1}{1+e^{-z_2^{(2)}}}, \frac{1}{1+e^{-z_3^{(2)}}} \right) \approx \begin{pmatrix} 1 & 0.948 & 0.962 & 0.599 \end{pmatrix},$ die 1 ist der Bias-Term. Schritt 4 - Gradientenberechnung (Backpropagation) $\frac{\partial}{\partial \Theta_{ij}^{(l-1)}} J(\Theta) = \frac{1}{n} \sum_{k=1}^n \delta_{ki}^{(l)} a_j^{(l-1)}$ $\delta_{11}^{(3)} \cdot a^{(2)} = -0.5 \cdot 0.948 = -0.474$ Schritt 5 Prüfen mittels Gradientenquotienten: $J(\Theta) = \frac{1}{2n} \left(\sum_{i=1}^n \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right)^2 \right)$ $= \frac{1}{2} (0.5 - 1)^2 = 0.125$ 			
Convolutional Neuronal Networks				
<p>Output Convolutional Filter: Output = $\sum_{i=0}^2 \sum_{j=0}^2$ Input $[i, j] \cdot$ Filter$[i, j]$ $= 0 \cdot 1 + 4 \cdot 2 + 2 \cdot 2 + 1 \cdot 13 \cdot 2 + 25 \cdot 4 + 5 \cdot 2 + 8 \cdot 1 + 6 \cdot 2 + 3 \cdot 1 = 169$</p> <p>Hauptidee: Extraktion lokaler Merkmale durch Faltungsoperationen.</p> <p>Faltungsschicht (Convolution Layer): wendet Filter (Kerne) an: $z = W * x + b$</p> <p>Pooling-Schicht: Reduktion der Dimensionalität (z.B. Max-Pooling oder Average-Pooling)</p> <p>Aktivierung: Typisch: ReLU $f(x) = \max(0, x)$</p> <p>Architektur-Beispiel: Input \rightarrow Conv \rightarrow ReLU \rightarrow Pool \rightarrow Dense \rightarrow Output</p> <p>Parameteranzahl: Abhängig von Filtergröße und Anzahl</p> <p>Vorteile:</p> <ul style="list-style-type: none"> Translation-Invarianz Weniger Parameter als vollständig verbundene Netze 				