

Unsupervised Deep Embedding for Clustering Analysis

Bachelorseminar Data Mining

Lukas Mahr

Ludwig-Maximilians-Universität München

1 Clustering of high dimensional data

2 Einleitung zu Neuronalen Netzen

- Idee
- Künstliches Neuron
- Layer/Schicht
- Aktivierungsfunktion
- Loss/Kostenfunktion
- Backpropagation mit Gradient descent

3 Autoencoders

- Idee
- Aufbau

4 Stecked Autoencoders

- Idee
- Aufbau

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

└ Roadmap

Roadmap

- Clustering of high dimensional data
- Einleitung zu Neuronalen Netzen
 - Idee
 - Künstliches Neuron
 - Layer/Schicht
 - Aktivierungsfunktion
 - Loss/Kostenfunktion
 - Backpropagation mit Gradient descent
- Autoencoders
 - Idee
 - Aufbau
- Stacked Autoencoders
 - Idee
 - Aufbau

Clustering of high dimensional data

■ Probleme

- unwichtige Features
- lange Cluster Zeiten
- Komplexität von z.B. KMeans
- $O(n^{dk+1})^{[1]}$ k =anz. Clusters, n =anz. Elemente, d =Dimension

■ Idee / Lösungsansatz

- Feature/Dimension Reduktion
- in Abhängigkeit der Clustere

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

└ Clustering of high dimensional data

└ Clustering of high dimensional data

viele Daten Punkte viele Distanzen zu berechnen schwierig zu visualisieren
ohne die Dimensionen zu reduzieren Komplexität von Kmeans die exponentiell ansteigt

Clustering of high dimensional data

- Probleme
 - unwichtige Features
 - lange Cluster Zeiten
 - Komplexität von z.B. KMeans
 - $O(n^{dk+1})^{[1]}$ k =anz. Clusters, n =anz. Elemente, d =Dimension
- Idee / Lösungsansatz
 - Feature/Dimension Reduktion
 - in Abhängigkeit der Clustere

Einleitung zu Neuronalen Netzen

Idee

— ? —

2022-02-02

- Unsupervised Deep Embedding for Clustering Analysis
 - └ Einleitung zu Neuronalen Netzen
 - └ Idee
 - └ Einleitung zu Neuronalen Netzen

Idee

— ? —

Künstlichen Neurons

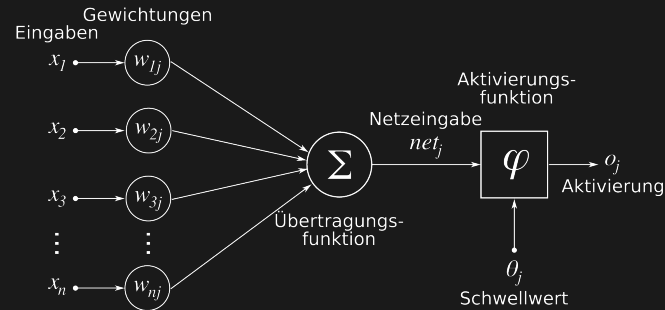
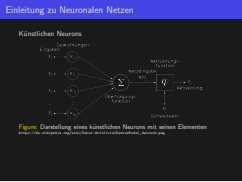


Figure: Darstellung eines künstlichen Neurons mit seinen Elementen
https://de.wikipedia.org/wiki/Datei:ArtificialNeuronModel_deutsch.png

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

- Einleitung zu Neuronalen Netzen
 - Künstliches Neuron
 - Einleitung zu Neuronalen Netzen



x_1, \dots, x_n sind die input variablen, jede der Eingabe variablen besitzt ein Gewicht, w_{1j}, \dots, w_{nj} . Diese werden Multipliziert und davon dann die summe berechnet. Hier die Übertragungsfunktion. Dazu wird ein Bias, in dem Fall der Schwellenwert gerechnet. Als letztes gibt es noch die Aktivierungsfunktion die meistens einen Wert zwischen 0 und 1 zurückgibt. Das ist dann der input für das nächste Neuron.

Layer/Schichten

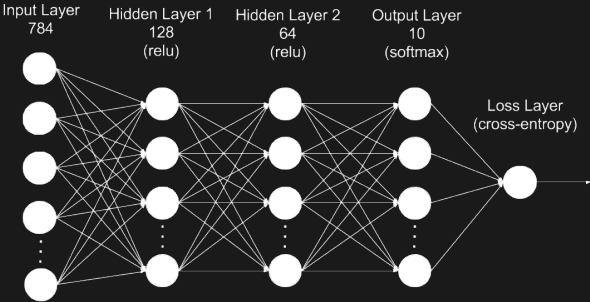
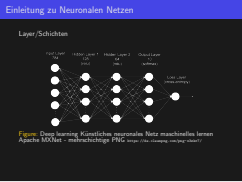


Figure: Deep learning Künstliches neuronales Netz maschinelles lernen
Apache MXNet - mehrschichtige PNG <https://de.cleapng.com/png-x3zkr7/>

2022-02-02

- Unsupervised Deep Embedding for Clustering Analysis
 - └ Einleitung zu Neuronalen Netzen
 - └ Layer/Schicht
 - └ Einleitung zu Neuronalen Netzen



Layer/Schicht sind mehrere Neuronen die mit allen Neuronen des nächsten Layer/Schicht verbunden sind. Alle Neuronen in einem Layer haben die gleiche Aktivierungsfunktion. Hidden Layer haben meistens die Aktivierungsfunktion rectified linear, da diese recht einfach und schnell zu berechnen ist. Das outputlayer hat meistens eine etwas kompliziertere Funktion wie softmax oder sigmoid. Abhängig von der Aufgabe des Netzwerkes. Letztes Layer hier direkt mit dem Loss

Aktivierungsfunktionen



Figure: Rectifier-Aktivierungsfunktion

https://de.wikipedia.org/wiki/Datei:Activation_rectified_linear.svg

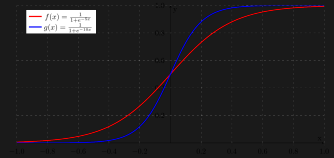


Figure: Sigmoide Funktion mit Steigungsmaß $a=5$ sowie $a = 10$

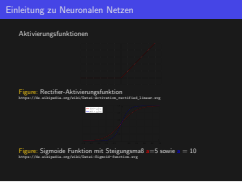
<https://de.wikipedia.org/wiki/Datei:Sigmoid-function.svg>

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

- └ Einleitung zu Neuronalen Netzen
 - └ Aktivierungsfunktion
 - └ Einleitung zu Neuronalen Netzen

alles negativ ist wird bei relu zu 0 während bei sigmoid, abhängig von der Steigung Werte zwischen -1 und 1 möglich sind



Loss/Kostenfunktion

Mean Squared Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

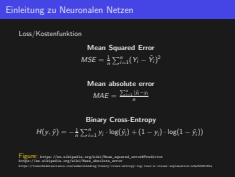
Mean absolute error

$$MAE = \frac{\sum_{i=1}^n |\hat{y}_i - y_i|}{n}$$

Binary Cross-Entropy

$$H(y, \hat{y}) = -\frac{1}{n} \sum_{i=1}^n y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)$$

Figure: https://en.wikipedia.org/wiki/Mean_squared_error#Predictor
https://en.wikipedia.org/wiki/Mean_absolute_error
<https://towardsdatascience.com/understanding-binary-cross-entropy-log-loss-a-visual-explanation-a3ac6025181a>



Man berechnet immer den unterschied zwischen den wahren labeln und den predicteden labeln um zu erkenne wie weit diese auseinander liegen. Es wird immer versucht den Loss zu minimieren. Also ein Minimum der Kostenfunktion zu finden. Die Parameter der Funktion, welche angepasst werden müssen sind alle weights und biases der einzelnen Neuronen und den Layern.

Backpropagation mit Gradient descent

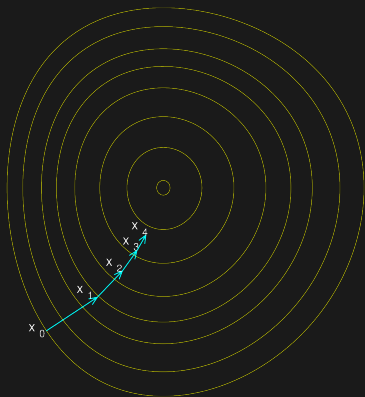


Figure: Illustration of gradient descent on a series of level sets
https://en.wikipedia.org/wiki/File:Gradient_descent.svg

Über Backpropagation wird hier mit z.B Gradient Descent die Loass funktion minimiert. Der Gradient der Loss/ Kostenfunktion wird für alle wights and biases gleichzeitig berechnet. Man kann sich das vorstellen, wie eine Kugel die man einen in einer Hügellandschaft rollen lässt ein kleinen schritten und zwischen den schritten immer nach der Steigung des Abhanges schaut und dabei versucht die Kugel in das tiefste Tal zu bekommen.

Autoencoders

Idee Aufbau Bottleneck

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

- └ Autoencoders
 - └ Aufbau
 - └ Autoencoders

Stacked Autoencoders

Idee Aufbau

2022-02-02

- Unsupervised Deep Embedding for Clustering Analysis
 - Stecked Autoencoders
 - Aufbau
 - Stacked Autoencoders

andere Clustering algorithmen ? andere
Dimensions-Reduktions-algorithmen

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

└─Stecked Autoencoders

└─Aufbau

└─Vorherige Arbeiten

Von wem ist das Paper

macvht hier kein sinn kommt am anfang

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

└─Stecked Autoencoders

└─Aufbau

└─Von wem ist das Paper

Von wem ist das Paper

macvht hier kein sinn kommt am anfang



k-means clustering

https://en.wikipedia.org/wiki/K-means_clustering#Complexity

2022-02-02

Unsupervised Deep Embedding for Clustering Analysis

└─ Referenzen

Referenzen

k-means clustering
https://en.wikipedia.org/wiki/K-means_clustering#Complexity