

Unsupervised Deep Embedding for Clustering Analysis

Bachelorseminar Data Mining

Lukas Mahr

Ludwig-Maximilians-Universität München

Roadmap

- 1 Clustern von Daten mit hohen Dimensionen
- 2 Einleitung zu Neuronalen Netzen
 - Idee
 - Künstliches Neuron
 - Layer/Schicht
 - Aktivierungsfunktion
 - Loss/Kostenfunktion
 - Backpropagation mit Gradient descent
- 3 Autoencoders
 - Idee
 - Aufbau
- 4 Unsupervised Deep Embedding for Clustering Analysis
 - Stacked Autoencoders
- 5 Deep embedded clustering
 - KMeans
 - Soft Assignments
 - Hilfsverteilung
 - Loss

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Roadmap

Roadmap

- Clustern von Daten mit hohen Dimensionen
- Einleitung zu Neuronalen Netzen
 - Idee
 - Künstliches Neuron
 - Layer/Schicht
 - Aktivierungsfunktion
 - Loss/Kostenfunktion
 - Backpropagation mit Gradient descent
- Autoencoders
 - Idee
 - Aufbau
- Unsupervised Deep Embedding for Clustering Analysis
 - Stacked Autoencoders
- Deep embedded clustering
 - KMeans
 - Soft Assignments
 - Hilfsverteilung
 - Loss

Clustern von Daten mit hohen Dimensionen

■ Problem

- Gaussian Mixture Models, KMeans
 - Abstandsmetriken sind beschränkt auf den ursprünglichen Datendimensionen
 - unwirksam wenn Datendimension hoch sind[1]
- Variationen von KMeans für Daten mit hohen Dimensionen
 - limitiert zu linearen Embeddings[2]
- Spectral clustering
 - Quadratische oder Super-quadratische Komplexität

■ Idee

- Neuronalesnetzwerk zum reduzieren der Dimensionen
 - nicht lineares mapping
- Clustern der reduzierten Daten
 - einfaches Clustering möglich, da Dimensionen reduziert
- Verbessern des NN und der Cluster durch Backpropagation

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Clustern von Daten mit hohen Dimensionen

└ Clustern von Daten mit hohen Dimensionen

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

Clustern von Daten mit hohen Dimensionen

- Problem
 - Gaussian Mixture Models, KMeans
 - Abstandsmetriken sind beschränkt auf den ursprünglichen Datendimensionen
 - unwirksam wenn Datendimension hoch sind[1]
 - Variationen von KMeans für Daten mit hohen Dimensionen
 - limitiert zu linearen Embeddings[2]
 - Spectral clustering
 - Quadratische oder Super-quadratische Komplexität
- Idee
 - Neuronalesnetzwerk zum reduzieren der Dimensionen
 - nicht lineares mapping
 - Clustern der reduzierten Daten
 - einfaches Clustering möglich, da Dimensionen reduziert
 - Verbessern des NN und der Cluster durch Backpropagation

Neuronale Netze

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Einleitung zu Neuronalen Netzen

└ Einleitung zu Neuronalen Netzen

Einleitung zu Neuronalen Netzen

Nicht-lineare statistische Modelle zur Informationsverarbeitung
Informationsverarbeitung umfasst hierbei unter anderem
 Klassifikation
 Prognosenerstellung
Units der Neuronalen Netze angelehnt an Neuronen
 Inputs zusammenfassen
 Mit Schwellenwert vergleichen bzw. aktivieren
Verbindungen zwischen Units angelehnt an Synapsen
 Gewichtung mit verstärkender oder schwächender Wirkung

2022-02-06

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

wofür braucht man neuronale Netze

- Klassifizierung von Daten
- Prognose eines bestimmten Wertes
- Neuronen an Neuronen in unserem Gehirn angelehnt
- Verbindungen zwischen Künstlichen Neuronen an Synapsen angelehnt
- supervised learning = überwachtes lernen
 - man besitzt Daten mit Beschriftungen (labels)
- unsupervised learning nicht überwachtes lernen
 - unbeschriftete Daten, also es sind keine labels vorhanden

Einleitung zu Neuronalen Netzen

Nicht-lineare statistische Modelle zur Informationsverarbeitung umfasst hierbei unter anderem
Klassifikation
Prognosenerstellung
Units der Neuronalen Netze angelehnt an Neuronen
Inputs zusammenfassen
Mit Schwellenwert vergleichen bzw. aktivieren
Verbindungen zwischen Units angelehnt an Synapsen
Gewichtung mit verstärkender oder schwächender Wirkung

Einleitung zu Neuronalen Netzen

Künstlichen Neurons

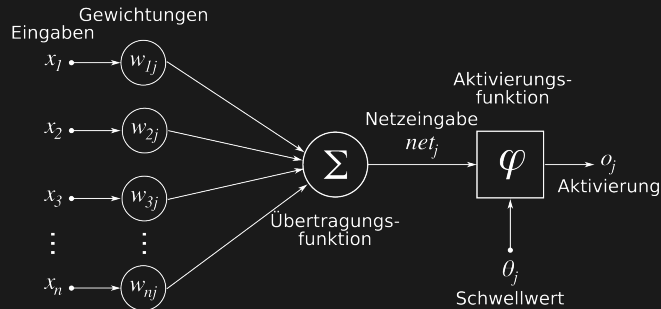
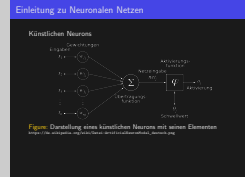


Figure: Darstellung eines künstlichen Neurons mit seinen Elementen
https://de.wikipedia.org/wiki/Datei:ArtificialNeuronModel_deutsch.png

2022-02-06

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

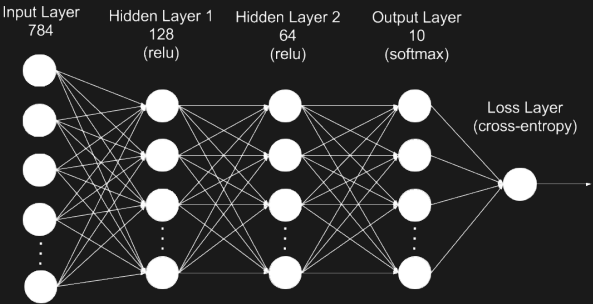
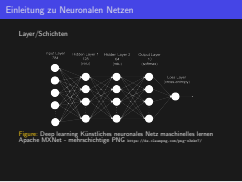


Figure: Deep learning Künstliches neuronales Netz maschinelles lernen
Apache MXNet - mehrschichtige PNG <https://de.cleapng.com/png-x3zkr7/>

2022-02-06

- 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

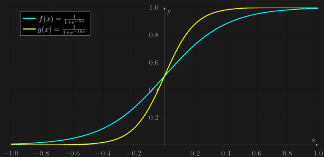
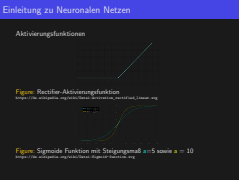


Figure: Sigmoide Funktion mit Steigungsmaß $a=5$ sowie $a = 10$
<https://de.wikipedia.org/wiki/Datei:Sigmoid-function.svg>

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

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

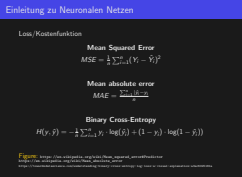
Binary Cross-Entropy

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

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>

Unsupervised Deep Embedding for Clustering Analysis

- └ Einleitung zu Neuronalen Netzen
 - └ Loss/Kostenfunktion
 - └ Einleitung zu Neuronalen Netzen



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

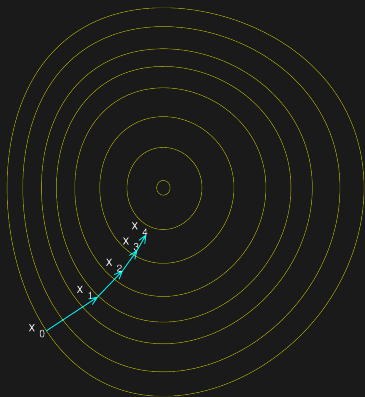


Figure: Illustration of gradient descent on a series of level sets
https://en.wikipedia.org/wiki/File:Gradient_descent.svg

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.

Autoencoder

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis
└ Autoencoders
 └ Idee
 └ Autoencoders

Autoencoders

Autoencoder

Ein Autoencoder ist ein feedforward Neural network, was versucht den input zu Kopieren. Das hört sich im ersten Moment nutzlos an, hat aber doch ein paar Anwendungsmöglichkeiten. Dazu gehört z.B Denosing oder Reduktion des Features Spaces

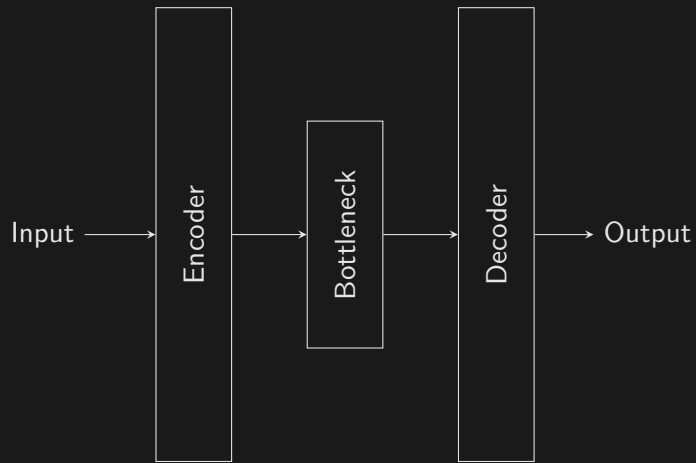


Figure: Einfaches Autoencoder Model

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Autoencoders

└ Aufbau

└ Autoencoders

Autoencoders

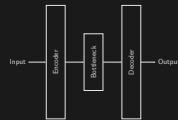


Figure: Einfaches Autoencoder Model

- Aufbau, Das Netzwerk besteht im Grunde aus 2 Teilen.
- Encoder und Decoder
- Encoder Transformiert die Input Daten in eine kleinere gewünschte Dimension
- Decoder Transformiert die Daten aus der Kleinen Dimension zurück in die Original Dimensionen.
- Hoffnung das der Encoder die Daten auf die Wichtigsten Features reduziert.

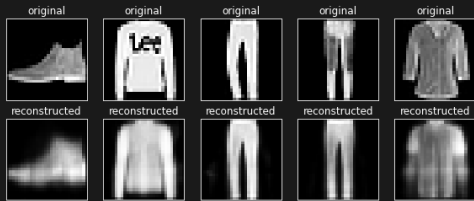


Figure: Original und Decoded Bilder

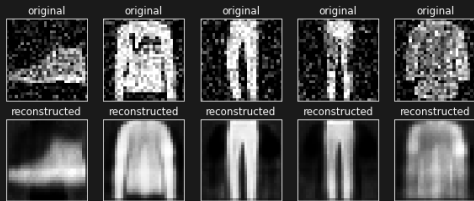


Figure: Noisy und Decoded Bilder

<https://github.com/Plutokezk/dec/blob/main/Autoencoders.ipynb>



Figure: Original und Decoded Bilder



Figure: Noisy und Decoded Bilder

Images: 11/12/13/14/15/16/17/18/19/20/21/22/23/24/25/26/27/28/29/30/31/32/33/34/35/36/37/38/39/40/41/42/43/44/45/46/47/48/49/50/51/52/53/54/55/56/57/58/59/60/61/62/63/64/65/66/67/68/69/70/71/72/73/74/75/76/77/78/79/80/81/82/83/84/85/86/87/88/89/90/91/92/93/94/95/96/97/98/99/100

- 784 input zu 10 zu 784 (28*28)
- Bottleneck size = 10 oder feature redutktion auf 10
- Random noise

DEC

Unsupervised Deep Embedding for Clustering Analysis

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis
└ Unsupervised Deep Embedding for Clustering Analysis
└ Unsupervised Deep Embedding for Clustering Analysis

DEC has two phases: (1) parameter initialization with a deep autoencoder parameter optimization (i.e., clustering), where we iterate between computing an auxiliary target distribution and minimizing the Kullback–Leibler (KL) divergence to it.

- Besteht aus 2 teilen
- embedding mapping in den kleineren raum Z
- Stacked Autoencoders
- clustering
- Clustern der embddeden Daten — Embeddings verbessern so wie die Cluster anpassen

Stacked Autoencoders

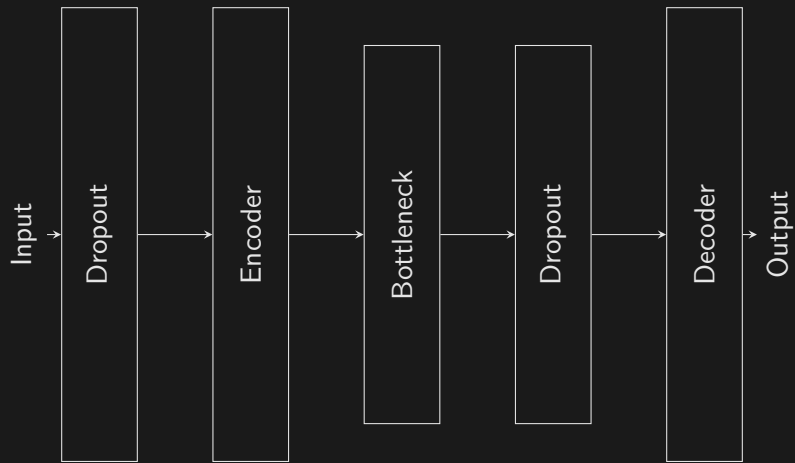
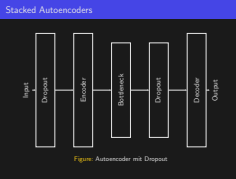


Figure: Autoencoder mit Dropout

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis
└ Unsupervised Deep Embedding for Clustering Analysis
 └ Stacked Autoencoders
 └ Stacked Autoencoders



Stacked Autoencoders ist genau das was der name sagt.
mehrere Autoencoder hintereinander.
man startet mit einem Autoencoder, den man trainiert.
dann entfernt man den decoder Teil
und setzt an diese stelle den nächsten Autoencoder, wird wieder trainiert
bis die gewünschte anzahl der stacks erreicht ist.
Besonderheit man trainiert die folgenden Autoencoder mit dem Output
des vorigen encoders.

Stacked Autoencoders

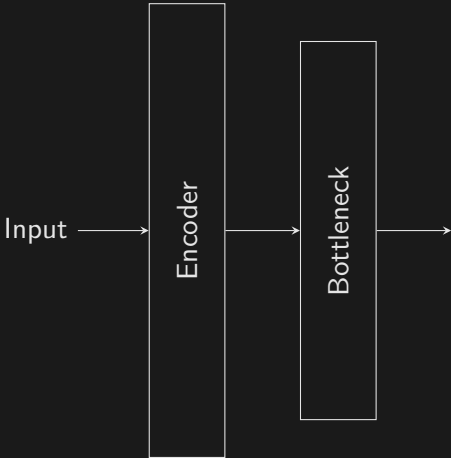


Figure: Encoder ohne Dropout

2022-02-06

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

Stacked Autoencoders

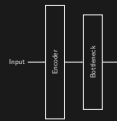


Figure: Encoder ohne Dropout

Encoder teil ohne noise/dropout und ohne dem decoder

Stacked Autoencoders

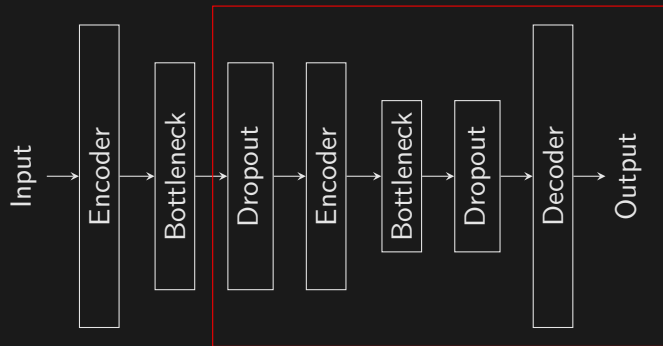
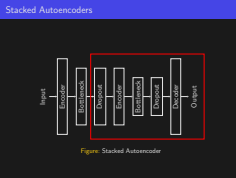


Figure: Stacked Autoencoder

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis
└ Unsupervised Deep Embedding for Clustering Analysis
└ Stacked Autoencoders
└ Stacked Autoencoders



Aktivierungsfunktion relu außer im letzten Decoder und Encoder Layer
sigmoid

für die zero-mean images

Zero-mean images - \bar{z} durchschnitt pro pixel über alle bilder im Datensatz
das bild wird dann von jedem bild abgezogen und damit liegt der durchschnitt bei den Bildern bei null deswegen zero mean images.

loss = least-square

nach dem Training des Autoencoders wird der Encoder als nicht lineares
mapping für die daten genutzt die man Clustern möchte

Deep embedded clustering

Deep embedded clustering

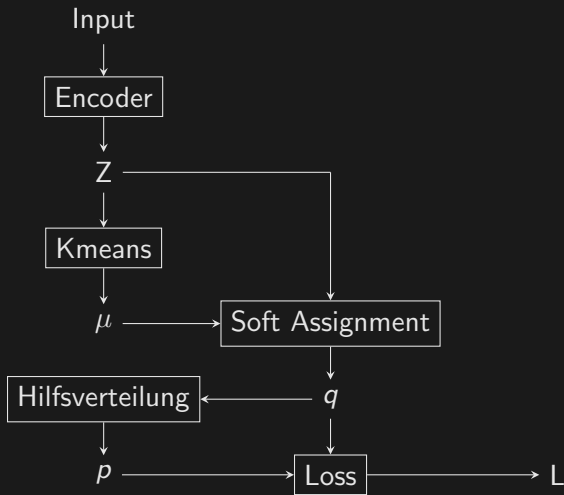


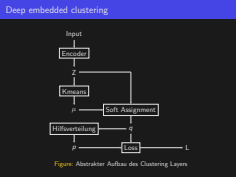
Figure: Abstrakter Aufbau des Clustering Layers

2022-02-06

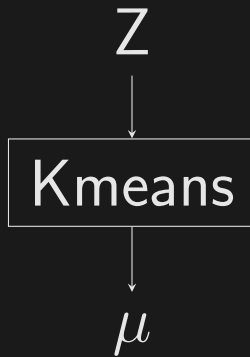
Unsupervised Deep Embedding for Clustering Analysis

└ Deep embedded clustering

└ Deep embedded clustering



- input daten durch sae werden auf kleinen feature raum Z abgebildet
- clustern der Daten Z es entstehen die Cluster Schwerpunkte μ
- soft assignments zwischen μ und z , wird zu q
- Hilfsverteilung aus q , wird zu p
- Loss zwischen q und p



2022-02-06

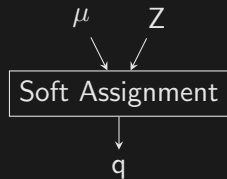
Unsupervised Deep Embedding for Clustering Analysis
└ Deep embedded clustering
 └ KMeans
 └ KMeans

KMeans



Clustern der Daten aus dem SAE

- Cluster Schwerpunkte für Soft Assignment
- eigentliche Cluster die man haben will



$$q_{ij} = \frac{(1 + ||z_i - \mu_j||^2 / \alpha)^{-\frac{\alpha+1}{2}}}{\sum_{j'} (1 + ||z_i - \mu'_j||^2 / \alpha)^{-\frac{\alpha+1}{2}}}$$

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

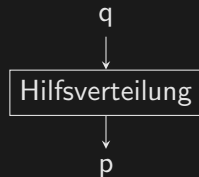
- └ Deep embedded clustering
 - └ Soft Assignments
 - └ Soft Assignments

Soft Assignments

$$q_{ij} = \frac{(1 + ||z_i - \mu_j||^2 / \alpha)^{-\frac{\alpha+1}{2}}}{\sum_{j'} (1 + ||z_i - \mu'_j||^2 / \alpha)^{-\frac{\alpha+1}{2}}}$$

Student t-Verteilung um die Ähnlichkeit zwischen z_i und μ_j zu messen
Freiheitsgrad ist 1[3]

q_{ij} ist die Wahrscheinlichkeit die z_i dem Cluster q_j zuzuordnen



$$p_{ij} = \frac{q_{ij}^2 / f_j}{\sum_{j'} q_{ij'}^2 / f_{j'}}$$

$$f_j = \sum_j q_{ij}$$

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Deep embedded clustering

└ Hilfsverteilung

└ Hilfsverteilung

Hilfsverteilung

```
graph TD; q --> HV[Hilfsverteilung]; HV --> p;
```

$$p_{ij} = \frac{q_{ij}^2 / f_j}{\sum_{j'} q_{ij'}^2 / f_{j'}}$$
$$f_j = \sum_j q_{ij}$$

Eigenschaften der Hilfsverteilung

- Prognosen Verstärken
 - mehr Gewicht auf genaue Datenpunkte
 - loss Verteilung von jedem Cluster Schwerpunkt normalisieren um zu verhindern das der versteckte Merkmals Bereich nicht verzehrt wird.
- man will das letzte Encoder Layer schützen.

f_j Soft-Cluster-Frequenzen

Erklärung wie die Eigenschaft erreicht werden im Fazit.

Kullback-Leibler-Divergenz



$$L = KL(P||Q) \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Deep embedded clustering

└ Loss

└ Loss

Loss

Kullback-Leibler-Divergenz



$$L = KL(P||Q) \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

Kullback-Leibler-Divergenz oder Kullback-Leibler-Abstand

- maß für Unterschiedlichkeit zweier Wahrscheinlichkeitsverteilungen
- hier als loss funktionen
- unterschied zwischen q und p
- der unterschied wird versucht zu minimieren.

$$\frac{\partial L}{\partial z_i} = \frac{\alpha+1}{\alpha} \sum_j (1 + \frac{\|z_i - \mu_j\|^2}{\alpha})^{-1} \times (p_{ij} - q_{ij})(z_i - \mu_j)$$

$$\frac{\partial L}{\partial \mu_j} = -\frac{\alpha+1}{\alpha} \sum_i (1 + \frac{\|z_i - \mu_j\|^2}{\alpha})^{-1} \times (p_{ij} - q_{ij})(z_i - \mu_j)$$

Stochastic Gradient Descent

$$\frac{\partial L}{\partial z_i} = \frac{\alpha+1}{\alpha} \sum_j (1 + \frac{\|z_i - \mu_j\|^2}{\alpha})^{-1} \times (p_{ij} - q_{ij})(z_i - \mu_j)$$

$$\frac{\partial L}{\partial \mu_j} = -\frac{\alpha+1}{\alpha} \sum_i (1 + \frac{\|z_i - \mu_j\|^2}{\alpha})^{-1} \times (p_{ij} - q_{ij})(z_i - \mu_j)$$

Gradient von L wird berechnet in Abhängigkeit von q_i und μ_j

$\frac{\partial L}{\partial z_i}$ wird dann für backpropagation genutzt.

wird Optimiert bis eine bestimmter Prozentsatz an Cluster punkten nicht mehr das Cluster wechselt, zwischen zwei fortlaufenden Iterationen.

andere Clustering algorithmen ? andere
Dimensions-Reduktions-algorithmen

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Deep embedded clustering

└ Optimierung

└ Vorherige Arbeiten

Von wem ist das Paper

Junyuan Xie
University of Washington

Ross Girshick
Facebook AI Research (FAIR)

Ali Farhadi
University of Washington

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

- └ Deep embedded clustering
 - └ Optimierung
 - └ Von wem ist das Paper

Von wem ist das Paper

Junyuan Xie
University of Washington

Ross Girshick
Facebook AI Research (FAIR)

Ali Farhadi
University of Washington



Steinbach, Michael, Ertöz, Levent, and Kumar, Vipin. The challenges of clustering high dimensional data. In *New Directions in Statistical Physics*, pp. 273–309. Springer, 2004.



Ye, Jieping, Zhao, Zheng, and Wu, Mingrui. Discriminative k-means for clustering. In *NIPS*, 2008.



van der Maaten, Laurens. Learning a parametric embedding by preserving local structure. In *International Conference on Artificial Intelligence and Statistics*, 2009.

2022-02-06

Unsupervised Deep Embedding for Clustering Analysis

└ Referenzen

└ Referenzen

- Steinbach, Michael, Ertöz, Levent, and Kumar, Vipin. The challenges of clustering high dimensional data. In *New Directions in Statistical Physics*, pp. 273–309. Springer, 2004.
- Ye, Jieping, Zhao, Zheng, and Wu, Mingrui. Discriminative k-means for clustering. In *NIPS*, 2008.
- van der Maaten, Laurens. Learning a parametric embedding by preserving local structure. In *International Conference on Artificial Intelligence and Statistics*, 2009.