

Neuronale Netze

Proseminar Data Mining

Lukas Krenz
Fakultät für Informatik
Technische Universität München
Email: lukas@krenz.land

Zusammenfassung—Das Paper ist eine kurze Einführung in Neuronale Netze, insbesondere in Multilayer Perceptrons inklusive Analyse ihrer praktischen Anwendung.

Schlüsselworte—Neuronale Netze, Multilayer Perceptrons, Backpropagation,

INHALTSVERZEICHNIS

I	Einleitung	1
II	Vom linearen Modell zum Neuronalen Netz	1
II-A	Das Perceptron als additives lineares Modell	1
II-B	Das Perceptron als einfaches ANN . .	1
II-C	Das XOR-Problem: nicht linear separierbare Probleme	1
II-D	Die Lösung: Neuronale Netze	2
III	Neuronale Netze	2
III-A	Die Aktivierungsfunktion	2
III-B	Die Kostenfunktion	2
III-C	Die Minimierung der Kostenfunktion - Backpropagation	2
III-D	Die numerische Methode - Gradient Descent und verwandte	2
IV	Neuronale Netze in der Praxis	3
IV-A	Problematik von lokalen Minima . . .	3
IV-B	Anzahl Hidden Layer	3
IV-C	Initialisierung der Gewichte	3
IV-D	Vermeidung von Overfitting	3
IV-E	Regression	3
IV-F	Klassifizierung	3
V	Zusammenfassung und Ausblick	3
Literatur		4

I. EINLEITUNG

Die Familie der Künstlichen Neuronalen Netze(ANN) ist wohl aktuell eins der spannendsten Forschungsgebiete im Bereich des maschinellen Lernens. In der letzten Zeit ist es die Methode, die in vielen Bereichen die besten Ergebnisse erzielt:

- Clustering

- Computer Vision
- Natural Language Processing.

Auch, wenn sie auf den ersten Eindruck kompliziert und unverständlich wirken, ist die mathematische und statistische Grundlage simpel.

Im folgenden wird der Weg vom linearen Modell zum Perceptron zum ANN gezeigt - mit jedem Schritt wird dabei die Komplexität, aber auch die Flexibilität erweitert. Nur die sogenannten Feed-Forward NNs werden hier beschrieben, wie sich später zeigen wird, sind diese auch mächtig genug, um viele Probleme zu lösen - auch wenn andere ANNs bessere Ergebnisse liefern können, sind die Grundlagen meistens übertragbar.

II. VOM LINEAREN MODELL ZUM NEURONALEN NETZ

Eines der klassischsten Modelle der Regressionsanalyse sind lineare Regression und die logistic Regression.

Es sind mathematisch simple Verfahren, die sogar analytisch lösbar sind - aber eben nicht immer ideale Ergebnisse erzielen: es gibt Probleme, die mit den Verfahren nicht lösbar sind.

Sehr simpel; lineares Modell als Input -> Output layer, dann lineare Funktion als Hidden Layer -> Lineares Modell

(Vergleich; Bild Logistic Regression vs Lineares Perceptron")

A. Das Perceptron als additives lineares Modell
Modell perceptron.

Es gibt viele Probleme, die von einem Perceptron nicht gelöst werden können. Es können nur Funktionen von linearen Modellen approximiert werden, die linear separierbar sind. Da das Perceptron als Verknüpfung von linearen Funktionen auch linear ist, kann auch es diese Klasse von Funktionen nicht lösen.

B. Das Perceptron als einfaches ANN
Verallgemeinerung zum ANN.

C. Das XOR-Problem: nicht linear separierbare Probleme
Theorem: Single Layer kann alle Funktionen annähern.

D. Die Lösung: Neuronale Netze

Neuronale Netze sind eine Verallgemeinerung von Perceptrons. Sie besitzen noch eine (oder mehrere) weitere Schicht(en) Neuronen, die jedoch nicht nur lineare Transformationen durchführen, sondern zusätzlich noch eine nicht lineare.

Diese Nicht-Linearität ermöglicht es, auch nicht linear separierbare Probleme zu lösen.

Durch die Nicht-Linearität und durch die komplexere Verknüpfung resultiert ein nicht lineares Optimierungsproblem, das meistens nicht mehr analytisch zu lösen ist: es werden numerische Verfahren benutzt.

III. NEURONALE NETZE

A. Die Aktivierungsfunktion

Die Funktion, die pro Knoten die nicht-lineare Transformation durchführt, heißt Aktivierungsfunktion. Häufig wird eine Funktion mit sigmoiden Erscheinungsbild gewählt, das heißt eine Funktion, die die Ergebnisse in ein bestimmtes Intervall "zusammenquetscht" (Quelle) und ein an ein "S" erinnerndes Erscheinungsbild hat. Die einfachste Funktion diesen Typs ist die so genannte logistische Funktion

$$\sigma_1(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

oder der hyperbolische Tangens

$$\sigma_2(x) = \tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}} = 2 * \sigma_1(2x) - 1 \quad (2)$$

Der Hyperbolische Tangens hat in der Praxis eine bessere Laufzeit (vgl. LeCun: Efficient Backprop)

Alternativ ReLu:

$$\sigma_3(x) = \max(0, x) \quad (3)$$

bis zu 6 mal schneller (TODO: Quelle)

...

B. Die Kostenfunktion

Wir benutzen eine Kostenfunktion, die den durchschnittlichen Fehler des ANN-Ergebnis quantifiziert.

Bei Regression wird meistens die Summe der quadrierten Terme benutzt (entspricht der euklidischen Norm):

$$J(\theta) = \sum_{k=1}^K \sum_{i=1}^N \left(\hat{Y}_i - Y_i \right)^2 \quad (4)$$

\hat{Y}_I entspricht dem Ergebnis des ANN, Y dem richtigen Ergebnis.

Bei Klassifizierungsproblemen wird üblicherweise die euklidische Norm, oder aber die negative Cross entropy benutzt [1]:

$$J(\theta) = - \sum_{k=1}^K \sum_{i=1}^N y_{ik} * \log f_k(x_i) \quad (5)$$

Bei Klassifikationsproblem ist Cross entropy fast immer die bessere Wahl. **To do (1)**

C. Die Minimierung der Kostenfunktion - Backpropagation

Die gewählte Kostenfunktion gilt es nun zu minimieren.

Backpropagation ist der bei Neuronalen Netzen benutzte Trainingsalgorithmus.

Er wird benutzt, um den Gradient der Kostenfunktion in Bezug auf alle Gewichte zu berechnen.

Dieser wird dann mit Hilfe einer geeigneten numerischen Methode benutzt, um die Gewichte zu optimieren.

Am Ende wird das Ergebnis der Optimierungsmethode benutzt, um die Gewichte zu aktualisieren.

To do (2)

D. Die numerische Methode - Gradient Descent und verwandte

Der Gradient wird oft mit der bekannten numerischen Methode Gradient Descent (eingedeutscht: Verfahren des steilen Abstiegs) benutzt.

Er ist relativ simpel:

To do (3)

Man unterscheidet beim Training von ANNs zwischen Batch und Stochastic Gradient Descent.

Bei der Batch Methode werden alle Gewichte auf einmal aktualisiert, bei Stochastic GD nur einzelne, zufällig ausgewählte Gewichte.

Das Gradientenverfahren hat das Problem, dass es oft in lokalen Minima steckenbleibt - das ist vor allem bei ANNs problematisch, weil die gewählte Kostenfunktion in den meisten Fällen mehrere lokale Minima besitzt. Auch wenn es nicht unbedingt notwendig, oder sogar erwünscht ist, einen globalen Extremwert zu finden (vgl. Overfitting), existieren bessere numerische Verfahren.

Eine Verbesserung zum gewöhnlichen Gradientenverfahren ist es, einen so genannten Momentumparameter zu benutzen, der auch vorherige Ergebnisse einbezieht, und somit das oft beobachtete oszillierende Verhalten des Gradientenverfahrens behebt.

$$\mathbf{b} = \mathbf{a} - \gamma \nabla F(\mathbf{a}) + \text{Momentum} \quad (6)$$

IV. NEURONALE NETZE IN DER PRAXIS

Weight Decay, Anfangswert

A. Problematik von lokalen Minima

Cost function hat mehrere Minima.

Lösung: Neu starten

Oft nicht globales Minimum erwünscht (siehe Overfitting)

B. Anzahl Hidden Layer

C. Initialisierung der Gewichte

Randomisiert - oft nicht gut

D. Vermeidung von Overfitting

Ein Problem, das in der Praxis auftreten kann, ist das so genannte Overfitting, das ist die übermäßige Anpassung an die Trainingsmenge. Oft ist das Modell zu stark angepasst und daher nicht mehr allgemein effizient.

Eine mögliche Lösung ist die L2-Regularisierung. Bei ihr wird an die Lossfunction ein weiterer Term angehängt:

Term einbauen + Quelle! (7)

.

Er skaliert mit der Summe aller Gewichte, es werden also komplexe Modelle bestraft
Bayesian View of L2-Reg. (wtf?)

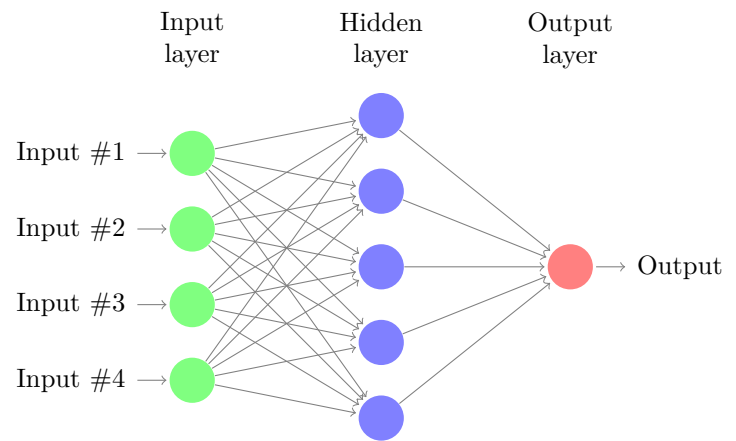
Eine weitere Lösung ist das so genannte early stopping, bei dem man nicht mit der Backpropagation aufhört, wenn ein (lokales) Minimum gefunden wurde, sondern dann, wenn die Performance des Modells bei dem Validierungsset optimal ist. Dabei wird oft die Iteration oft beendet, wenn bei der Fehlerfunktion gar kein Minimum vorliegt.

E. Regression

F. Klassifizierung

Cross entropy.

Soft-Max-F unction



$$\begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \quad \begin{pmatrix} h_{1,1} \\ h_{1,2} \\ h_{1,3} \\ \vdots \\ h_{1,m} \end{pmatrix} \quad \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

V. ZUSAMMENFASSUNG UND AUSBLICK

blabla

TO DO...

- ☐ 1 (p. 2): Warum ist Cross entropy besser?
- ☐ 2 (p. 2): Backpropagation
- ☐ 3 (p. 2): Gradient Descent (primär Formel)

LITERATUR

- [1] T. Hastie, R. Tibshirani, J. Friedman, T. Hastie, J. Friedman, and R. Tibshirani, *The elements of statistical learning*. Springer, 2009, vol. 2, no. 1.