

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,

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

A. Das Perceptron als additives lineares Modell

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.

Eine Weiterentwicklung ist das so genannte Perceptron. Es besteht aus einem Eingabe-Vektor x , einem Gewichtsvektor w und einem Bias-Wert b . Es ist ein binärer Klassifizierer, der das Ergebniss mit Hilfe der folgenden Formel berechnet:

$$y = \begin{cases} 1 & w * x + b > 0 \\ 0 & \text{ansonsten} \end{cases} \quad (1)$$

To do (1)

B. Das Perceptron als einfaches ANN

Wie wir später sehen werden, ist das Perceptron eine einfache Variante eines Feed-Forward-Neuronalen-Netzes mit der so genannten Heaviside-Funktion

$$\sigma = \begin{cases} 1 & n \geq 0 \\ 0 & n \leq 0 \end{cases} \quad (2)$$

als Aktivierungsfunktion.

Auf den benutzen Trainingsalgorithmus werden wir nicht eingehen; der später betrachtete Backpropagation-Algorithmus für allgemeine ANNs ist auf Perceptrons übertragbar.

C. Das XOR-Problem: nicht linear separierbare Probleme

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

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. Aufbau und Namenskonventionen

Ein Neuronales Netz besteht aus mehreren Schichten Neuronen. Einem Input Layer, einen (oder mehreren) Hidden Layer, und einem Output Layer. Wenn man von

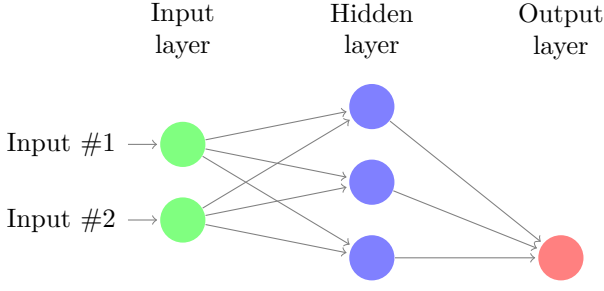


Abbildung 1. Ein 2-schichtiges MLP.

einem k-schichtigen ANN spricht, bezeichnet k die Anzahl der Schichten, ohne Input-Layer.

Sowohl Input, als auch Hidden-Layer bestehen aus Neuronen mit einer Aktivierungsfunktion, die die Eingangsdaten transformiert.

Jeder Layer ist mit den anderen Elementweise verbunden. Jede Verbindung hat ein so genanntes Gewicht. Daten fließen nur von links nach rechts.

Wir benutzen folgende Namenskonvention: m ist die Anzahl der Parameter, $x \in \mathbb{R}^{m \times 1}$ ist der Input-Vektor, $w_i \in \mathbb{R}^{\text{Anzahl Neuronen} \times \text{Anzahl Neuronen Layer davor}}$ der Gewichtsvektor, der die i -te Schicht mit der $i + 1$ Schicht verbindet; $b \in \mathbb{R}^{\text{Anzahl Neuronen} \times 1}$ der Bias-Vektor; \hat{y} ist der Output des Netzwerkes, y der Referenzwert (Ergebnis im Trainingsset).

B. Feedforward

$$\text{net}_j = \sum_{i=1}^d x_i w_{j,i} + w_{j,0} = \sum_{i=0}^d x_i w_{j,i} = w_j^t x \quad (3)$$

$$y_j = \sigma(\text{net}_j) \quad (4)$$

[1]

C. 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 dieses Typs ist die so genannte logistische Funktion

$$\sigma_1(x) = \frac{1}{1 + e^{-x}} \quad (5)$$

oder der hyperbolische Tangens

$$\sigma_2(x) = \tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}} = 2 * \sigma_1(2x) - 1 \quad (6)$$

Die logistische Funktion ist nicht 0-zentriert. Das hat den Effekt, dass, falls alle Eingangswerte des Neurons, der

Gradient entweder komplett positiv, oder negativ ist. Das führt unter Umständen zu unerwünschten Oszillationen bei der Optimierung. Der hyperbolische Tangens ist 0-zentriert, hat aber einen Nachteil mit der logistischen Funktion gemeinsam: sehr kleine oder sehr große Werte haben einen sehr kleinen Gradienten zur Folge; dies führt zum "Absterben" des Neurons.

Eine alternative Aktivierungsfunktion ist die Rectifier-Funktion

$$\sigma_3(x) = \max(0, x). \quad (7)$$

Glorot et al argumentieren, dass die logistische Funktion zwar biologisch plausibler ist als der hyperbolische Tangens, aber die Rectifier Funktion noch näher an der Funktionsweise biologischer Neuronen ist. Was für uns relevanter ist: sie ist effizienter zu berechnen; also in der Praxis oft eine bessere Wahl. [2]

D. 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):

$$J(\theta) = \sum_{k=1}^K \sum_{i=1}^N (\hat{Y}_i - Y_i)^2 \quad (8)$$

\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 [3]:

$$J(\theta) = - \sum_{k=1}^K \sum_{i=1}^N y_{ik} * \log f_k(x_i) \quad (9)$$

Bei Klassifikationsproblem ist Cross entropy fast immer die bessere Wahl. **To do** (2)

E. 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. [4]

Alle relevanten Fehlerfunktionen haben die Form

$$J = \sum_n J_n \quad (10)$$

das heißt, der Gesamtfehler ist die Summe aller Einzelfehler. Wir nehmen weiterhin an, dass der Fehler als ableitbare Funktion der Output-Werte darstellbar ist.

$$J_n = J_n(y_1, \dots, y_c) \quad (11)$$

Durch mehrfaches Anwenden der Kettenregel bekommen wir die so genannte Backpropagation-Formel:

$$\theta_j = \sigma'(a_j) \sum_k w_{k,j} \sigma_k \quad (12)$$

To do (3)

[5]

F. 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 (4)

$$x_{n+1} = x_n - \gamma_n \nabla F(x_n) \quad (13)$$

Auf ANNs angewandt, sehen die Formeln dann folgendermaßen aus:

$$\Delta w_{j,i} = -\eta \delta_j x_i \quad (14)$$

$$\Delta w_{j,i} = -\eta \sum_n \delta_j^n x_i^n \quad (15)$$

[5]

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.

$$\mathbf{b} = \mathbf{a} - \gamma \nabla F(\mathbf{a}) + \text{Momentum} \quad (16)$$

IV. NEURONALE NETZE IN DER PRAXIS

Weight Decay, Anfangswert

A. Problematik von lokalen Minima

Die Kostenfunktion hat oft mehrere lokale Minima - das Gradientenverfahren kann dann in einem solchen steckenbleiben, das heißt, es konvergiert unter Umständen nicht zum gesuchten globalen Minimum. Auch wenn das globale Minimum nicht immer erwünscht ist (vgl. early stopping), ist es unter Umständen doch besser, es zu finden.

Die einfachste Lösung wäre natürlich, ein anderes numerisches Verfahren zu benutzen. Doch andere Verfahren sind meistens komplizierter zu implementieren, und unter Umständen sogar langsamer als SGD. Wenn man beim SGD in einem Minimum stecken bleibt, bietet es sich an, einfach mit neu initialisierten Gewichtsparametern neuzustarten.

B. Anzahl Hidden Layer

C. Initialisierung der Gewichte

Üblicherweise werden die Gewichte bei Start mit kleinen Werten nahe bei 0 initialisiert - ein Startwert von 0 würde bedeuten, dass keine Daten durch die Verbindung fließen. Es gibt aber oft gute Heuristiken für Initialisierungen, die sich in der Praxis bewährt haben. So haben zum Beispiel ...beobachtet, dass bei der logistischen Funktion besonders gut π geeignet ist, bei der tanh oft $6 \times \pi$.

D. Vermeidung von Overfitting

Ein Problem, dass 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 L1/L2-Regularisierung. Bei ihr wird an die Kostenfunktion ein weiterer Term angehängt:

$$J_{L1} = J + \lambda \sum_i \theta_i \quad (17)$$

$$J_{L2} = J + \lambda \sum_i \theta_i^2 \quad (18)$$

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

Aus der Bayesianischen Perspektive betrachtet, ist eine L2-Regulierung gleichbedeutend damit, dass wir eine prior-Distribution annehmen.

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

V. ZUSAMMENFASSUNG UND AUSBLICK

Wir haben nur MLPs betrachtet, in neueren Forschungsergebnissen werden oft ANNs mit anderer Topologie benutzt, das heißt, die Daten fließen nur in eine Richtung. Relevant sind aber auch andere Netzwerkarten wie RNNs, RBM, etc. Sie besitzen einen komplizierteren Aufbau; aber im Grunde sind die meisten genannten Grundlagen auch für diese Arten von Netzwerken relevant. Die Zukunft bleibt spannend, ANNs sind sehr vielversprechend.

TO DO...

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

LITERATUR

- [1] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern classification*. John Wiley & Sons, 2012.
- [2] X. Glorot, A. Bordes, and Y. Bengio, “Deep sparse rectifier networks,” in *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics. JMLR W&CP Volume*, vol. 15, 2011, pp. 315–323.
- [3] T. Hastie, R. Tibshirani, J. Friedman, T. Hastie, J. Friedman, and R. Tibshirani, *The elements of statistical learning*. Springer, 2009, vol. 2, no. 1.
- [4] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Cognitive modeling*, vol. 5, 1988.
- [5] C. M. Bishop, *Neural networks for pattern recognition*. Clarendon press Oxford, 1995.