



HUMBOLDT UNIVERSITÄT ZU BERLIN

HAUSARBEIT

**Vorhersage von
Ausfallwahrscheinlichkeiten mit
neuronalen Netzen**

Thomas Siskos (580726)

DATENANALYSE II

Dozent:

Dr. Sigmund KLINKE

28. März 2018

Inhaltsverzeichnis

1	Einleitung	2
2	Die Creditreform Datenbank	4
3	Methoden	7
3.1	Neuronale Netze	7
3.2	Architekturen und Hyperparameter	9
3.3	Ergebnisse	9
4	Zusammenfassung	9

Tabellenverzeichnis

1	Variablen des Creditreform Datensatzes	3
2	Definitionen der finanziellen Kennzahlen	5

Abbildungsverzeichnis

1	Neuronales Netz mit einer versteckten Schicht.	8
2	Neuronales Netz mit zwei versteckten Schichten.	10
3	Neuronales Netz mit fünf versteckten Schichten.	10

1 Einleitung

Die Ausfallwahrscheinlichkeit ist ein Begriff der Finanzen und beschreibt die Wahrscheinlichkeit, dass ein Kreditnehmer, innerhalb eines Zeitrahmens, nicht in der Lage sein wird seine Verpflichtungen zum vorher festgelegten Termin einzuhalten. Die möglichst genaue Schätzung und Vorhersage eines solchen Ausfalls ist von entscheidender Bedeutung. Die Bepreisung von Vermögenswerten, das Einschätzen der Risiken von Krediten und Kreditportfolios sowie die Wertschätzung anderer Finanz-Produkte hängen maßgeblich von der Präzision der geschätzten Ausfallwahrscheinlichkeiten ab (Miao et al., 2008). Es gibt zwei wesentliche Denkschulen bei der Analyse von Kreditausfällen, der markt-basierte und der statistische Ansatz. Die markt-basierte Art modelliert Kreditausfälle mithilfe struktureller Modelle, wohingegen der statistische Ansatz empirische Methoden bemüht, historische Daten auszuwerten. Diese Daten können beispielsweise aus der Buchhaltung, beziehungsweise aus den Bilanzen von Unternehmen gewonnen werden (Härdle et al., 2012).

Zur Quantifizierung der Ausfallwahrscheinlichkeit werden in der Literatur oft Kennzahlen verwendet, die vornehmlich eine oder mehrere Bilanzpositionen gegen eine oder mehrere andere ins Verhältnis setzen. Diese Kennzahlen haben sich in vergangenen Studien oft als hilfreich erwiesen (Altman, 1968). In dieser Arbeit werden wir die 28 finanziellen Kennzahlen bestimmen, die von Zhang und Härdle (2010) verwendet wurden. Wir werden mithilfe dieses transformierten, des untransformierten sowie eines zusammengetzten Datensatzes neuronale Netze unterschiedlicher Architekturen trainieren und miteinander vergleichen. Insbesondere verwenden wir klassische neuronale Netze mit einer versteckten Schicht, Netze mit zwei versteckten Schichten und Netze mit fünf versteckten Schichten. Die verschiedenen Architekturen wurden mithilfe der Programmiersprache `Python`, genauer mit dem Modul `tensorflow`, erzeugt. Der Quellcode für die neuronalen Netze ist auf github.com/thsis/DAII einsehbar.

Der folgende Abschnitt beschreibt den Datensatz der Creditreform Datenbank, sowie die Art und Weise wie die Daten bereinigt und transformiert wurden um die 28 Finanz-Kennzahlen zu bestimmen. Abschnitt 3 erläutert die Theorie und Anwendung neuronaler Netze auf die transformierten und untransformierten Daten. Der letzte Abschnitt enthält eine abschließende Zusammenfassung und kritische Würdigung der Ergebnisse.

Tabelle 1: Variablen des Creditreform Datensatzes

Variable	Bedeutung
ID	Kennnummer jedes Unternehmens
T2	Solvenz-Status (solvent:0, insolvent:1)
JAHR	Jahr
VAR1	Scheck, Kassenbestand
VAR2	Vorräte - Gesamt
VAR3	Umlaufvermögen
VAR4	Sachanlagen - Gesamt
VAR5	Immaterielle Vermögensbestände
VAR6	Gesamtvermögen
VAR7	Forderungen aus Lieferung und Leistung
VAR29	Forderungen ggü. Unternehmen mit Beteiligungsverhältnis
VAR8	Grundstücke und Bauten
VAR9	Eigenkapital
VAR10	Gesellschafterdarlehen
VAR11	Pensionsrückstellungen
VAR12	kurzfristige Verbindlichkeiten - Gesamt
VAR13	langfristige Verbindlichkeiten - Gesamt
VAR14	Bankschulden
VAR15	Verbindlichkeiten aus Lieferung und Leistung
VAR30	Verbindlichkeiten ggü. Unternehmen mit Beteiligungsverhältnis
VAR16	Umsätze
VAR17	Vertriebs- / Verwaltungsaufwand
VAR18	Abschreibungen
VAR19	Zinsaufwendungen
VAR20	Gewinn vor Zins und Steuern (EBIT)
VAR21	Betriebsgewinn
VAR22	Jahresüberschuss
VAR23	(Lager-) Bestandsveränderungen
VAR24	Veränderungen der Verbindlichkeiten ggü. Vorjahr
VAR25	Veränderungen Bargeld/Kassenbestand/flüssige Mittel
VAR26	Branchenzugehörigkeit
VAR27	Rechtsform
VAR28	Anzahl Mitarbeiter

2 Die Creditreform Datenbank

Die Creditreform Datenbank enthält Daten für 20.000 solvente und 1.000 insolvente deutsche Firmen aus den Jahren 1997 bis 2007. Der Datensatz wurde durch das Labor für empirische und quantitative Forschung (LEQR) der Humboldt Universität zu Berlin bereitgestellt. Die enthaltenen Variablen stammen aus den Bilanzen der Unternehmen und stellen für potentielle Investoren die Hauptgrundlage für Analysen dar. Ein Unternehmen wird entweder mit sich selbst verglichen, indem der zeitliche Verlauf der Bilanzposten untersucht wird oder das Unternehmen wird mit ähnlichen Firmen verglichen indem eine Auswahl finanzieller Kennzahlen betrachtet wird (Berk & DeMarzo, 2016).

Rund die Hälfte der Daten stammt aus den Jahren 2001 und 2002. Da 1996 keine insolventen Unternehmen vorliegen, werden alle Beobachtungen dieses Jahres gelöscht. Der Großteil der verbleibenden Unternehmen ist entweder im Baugewerbe, im Handel, in der Industrie oder im Immobilienwesen tätig. Andere Kategorien umfassen beispielsweise Branchen wie Landwirtschaft und Bergbau, Elektrizität-, Gas- und Wasserversorgung, die Gastronomie, Logistik und soziale Dienstleistungen. Alle Unternehmen die zu diesen Kategorien gehören werden von der folgenden Analyse ausgeschlossen, um die Schichtung des Trainingsdatensatzes nicht unnötig zu komplizieren. Außerdem werden sowohl die kleinsten, als auch die größten Unternehmen entfernt. Betrachtet werden nur Unternehmen, deren Gesamtvermögen zwischen 10.000 und 10.000.000 liegen. Die kleinsten Unternehmen werden entfernt, da deren finanzielle Lage oft von den Finanzen einer einzelnen verantwortlichen Person, typischerweise die Eigentümerin oder der Eigentümer, abhängt. Die größten Unternehmen werden hingegen entfernt, da sie in Deutschland nur in den allerseltensten Fällen Gefahr laufen in die Zahlungsunfähigkeit zu geraten. Des Weiteren werden Unternehmen entfernt, bei denen während der Berechnung der finanziellen Kennzahlen Nullen im Nenner auftreten (Chen, 2010).

Die verbleibenden Unternehmen gliedern sich in verschiedene Sektoren, von denen die vier häufigsten im Folgenden analysiert werden. Von den 9567 solventen Unternehmen sind 35,9% in der Industrie, 34,1% im Handel, 19,5% im Baugewerbe und 10,4% in der Immobilienbranche tätig. Von den 782 insolventen Unternehmen sind 45,0% im Baugewerbe, 28,2% in der Industrie, 21,6% im Handel und 5,1% in der Immobilienbranche tätig.

Tabelle 2: Definitionen der finanziellen Kennzahlen

Name	Formel	Kennzahl
x1	$\text{VAR22}/\text{VAR6}$	Gesamtkapitalrentabilität (ROA)
x2	$\text{VAR22}/\text{VAR16}$	Nettogewinnmarge
x3	$\text{VAR21}/\text{VAR6}$	Betriebsgewinnmarge
x4	$\text{VAR21}/\text{VAR16}$	
x5	$\text{VAR20}/\text{VAR6}$	EBITDA
x6	$(\text{VAR20}+\text{VAR18})/\text{VAR6}$	
x7	$\text{VAR20}/\text{VAR16}$	Eigenmittelquote (einfach) Eigenmittelquote (angepasst)
x8	$\text{VAR9}/\text{VAR6}$	
x9	$(\text{VAR9}-\text{VAR5})/(\text{VAR6}-\text{VAR5}-\text{VAR1}-\text{VAR8})$	Nettoverschuldung
x10	$\text{VAR12}/\text{VAR6}$	
x11	$(\text{VAR12}-\text{VAR1})/\text{VAR6}$	Schuldenquote
x12	$(\text{VAR12}+\text{VAR13})/\text{VAR6}$	
x13	$\text{VAR14}/\text{VAR6}$	Zinsdeckungsgrad
x14	$\text{VAR20}/\text{VAR19}$	
x15	$\text{VAR1}/\text{VAR6}$	Liquiditätsgrad
x16	$\text{VAR1}/\text{VAR12}$	
x17	$(\text{VAR3}-\text{VAR2})/\text{VAR12}$	Quick Ratio Current Ratio
x18	$\text{VAR3}/\text{VAR12}$	
x19	$(\text{VAR3}-\text{VAR12})/\text{VAR6}$	Kapitalumschlag
x20	$\text{VAR12}/(\text{VAR12}+\text{VAR13})$	
x21	$\text{VAR6}/\text{VAR16}$	Lagerumschlag
x22	$\text{VAR2}/\text{VAR16}$	
x23	$\text{VAR7}/\text{VAR16}$	Forderungsumschlag
x24	$\text{VAR15}/\text{VAR16}$	
x25	$\log(\text{VAR6})$	Verbindlichkeitenumschlag
x26	$\text{VAR23}/\text{VAR2}$	
x27	$\text{VAR24}/(\text{VAR12}+\text{VAR13})$	Proxy für die Unternehmensgröße Gestaffelte Prozentuale Lagerbestandsänderung Gestaffelte Prozentuale Forderungsänderung Gestaffelte Prozentuale Änderung des Cash Flow
x28	$\text{VAR25}/\text{VAR1}$	

Anschließend werden mithilfe der verbleibenden Unternehmen die finanziellen Kennzahlen ermittelt. Die Variablen **x1-x7** gehören zu den sogenannten Rentabilitätsverhältnissen. Rentabilitätsverhältnisse haben sich in der Vergangenheit als besonders starke Prädiktoren für Kreditausfälle erwiesen. Zum Beispiel gewährt die Gesamtkapitalrentabilität (return on assets, ROA) **x1** einen Einblick in die Umsatzstärke eines Unternehmens im Vergleich zu dessen Kosten. So signalisiert ein höherer Wert der Kennzahl, dass ein Unternehmen in der Lage ist mehr Geld mit weniger Mitteln zu verdienen. Die Nettogewinnmarge **x2** hingegen veranschaulicht den Anteil des Umsatzes, den das Unternehmen als Einnahmen einbehält. Ein hoher Wert geht mit einem profitablen Unternehmen einher, das seine Kosten zu kontrollieren versteht (Chen et al., 2011).

Eine weitere Reihe der Kennzahlen beschreibt die sogenannte Hebelwirkung. Damit ist das Ausmaß gemeint, in dem ein Unternehmen auf Schulden als Finanzierungsquelle angewiesen ist (Berk & DeMarzo, 2016). Da Unternehmen Schulden und Eigenkapital kombinieren, um ihre Aktivitäten zu finanzieren, erweisen sich Kennzahlen über ebenjene Hebelwirkung als hilfreiche Werkzeuge um die Zahlungsfähigkeit eines Unternehmens einzuschätzen. Zu den Kennzahlen der Hebelfinanzierung gehören die Variablen **x8-x14**. Beispielsweise misst die Nettoverschuldung **x11** die Höhe der

kurzfristigen Verpflichtungen, welche nicht durch die liquiden Vermögensbestände gedeckt sind, als prozentualer Anteil des Gesamtvermögens. Somit misst diese Kennzahl auch die kurzfristige Liquidität eines Unternehmens (Chen et al., 2011).

Die sechs Folgenden Verhältnisse **x15-x20** gehören in den Bereich der Liquiditäts-Kennzahlen. Liquidität ist eine weit verbreitete Variable, die in vielen Kreditentscheidungen eine wichtige Rolle spielt. Liquidität beschreibt die Möglichkeiten eines Unternehmens Vermögensbestände in kurzer Zeit in Bargeld umzuwandeln. Die wohl wichtigste Kennzahl für die Liquidität ist der Anteil der Kassenbestände am Gesamtvermögen **x15**. Ein weiterer wichtiger Gradmesser für die Liquidität eines Unternehmens ist der sogenannte Quick-Ratio **x17**. Mithilfe des Quick-Ratio versucht man einzuschätzen, ob ein Unternehmen über ausreichend liquide Mittel verfügt um insbesondere kurzfristige Zahlungen zu decken. Ein hoher Quick-Ratio indiziert, dass es für das Unternehmen unwahrscheinlich ist, kurzfristig in Zahlungsnot zu geraten (Berk & DeMarzo, 2016).

Einen weiteren wichtigen Typus von betriebswirtschaftlichen Kennzahlen stellen die Aktivitätskennzahlen **x21-x24** dar. Aktivitätskennzahlen messen die Effizienz mit der ein Unternehmen eigene Ressourcen aufwendet, um Umsatz mithilfe seiner Vermögensbestände zu generieren (Chen et al., 2011).

Zusätzlich berechnen wir den Logarithmus des Gesamtvermögens **x25**. Dieser Risikoindikator stellt die Größe eines Unternehmens dar und versetzt uns in die Lage große, mittlere und kleine Unternehmen miteinander in Beziehung zu setzen. Als letzte Gruppe betriebswirtschaftlicher Kennzahlen berechnen wir die gestaffelten prozentualen Änderungen des Cash-Flow, des Lagerbestandes und der Forderungen im Vergleich zum Vorjahr **x26-x28**. (Chen et al., 2011).

Um Einflüsse von Ausreißern auf die neuronalen Netze zu eliminieren, werden extreme Werte für die verschiedenen Verhältnisse durch das 0.05- bzw. das 0.95-Quantil ersetzt. Präziser ausgedrückt, folgen wir der Regel, wenn $x_{ij} < q_{0.05}(x_j)$, dann setze $x_{ij} \stackrel{!}{=} q_{0.05}(x_j)$. Beziehungsweise, wenn $x_{ij} > q_{0.95}(x_j)$, dann setze $x_{ij} \stackrel{!}{=} q_{0.95}(x_j)$. Wobei $x_i, i \in \{1, \dots, N\}$ für jeden einzelnen Wert einer Kennzahl x_j und $q_k(x_j), j \in \{1, \dots, 28\}, k = 0.05, 0.95$ für die jeweiligen Quantile der Kennzahl x_j des Datensatzes steht.

3 Methoden

3.1 Neuronale Netze

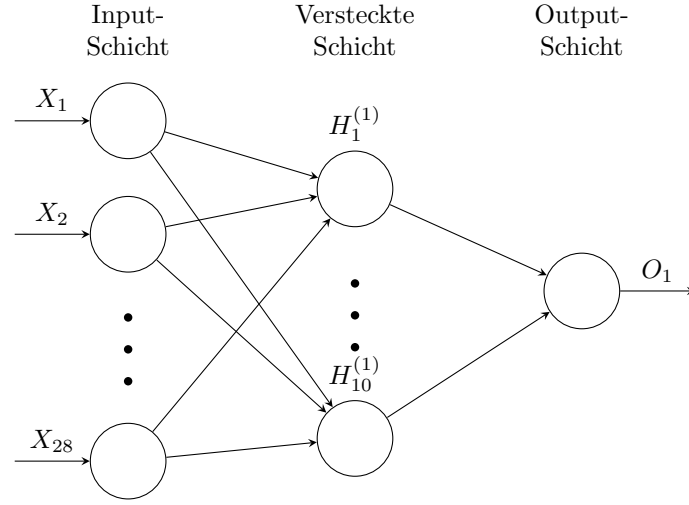
Künstliche neuronale Netze versuchen die Struktur von Gehirnen nachzubilden. Stark vereinfacht enthält ein jedes Gehirn sogenannte Neuronen, welche über zwei Zustände verfügen können. Ein Neuron ist entweder aktiviert oder nicht. In einem Gehirn verändern Neuronen ihren Zustand als Reaktion auf einen chemischen oder elektrischen Stimulus. Das Netz, das die Neuronen innerhalb des Gehirns eines Menschen bilden ist ein enormes Gespinnst, in dem der Zustand eines Neurons das Ergebnis tausender anderer Neuronen sein kann. Führt ein Stimulus dazu, dass bestimmte Verbindungen wiederholt aktiviert werden, verfestigt sich deren Verbindung. Das führt dazu, dass bei einem ähnlichen Stimulus ebenfalls dieselben Verbindungen aktiviert werden und zu demselben Outputzustand führen. Dieses Verhalten nennen wir Lernen.

Künstliche neuronale Netze vereinfachen die Vorgänge des Gehirns stark. Ein künstliches Neuron wird durch eine Aktivierungsfunktion simuliert, deren Bildmenge das Verhalten eines Schalters emuliert. Wir verlangen von der Aktivierungsfunktion, dass sie über mindestens zwei deutlich verschiedene Zustände verfügt. Typischerweise ist der Output einer Aktivierungsfunktion zum Beispiel entweder Null oder Eins, Null oder größer Null, Minus Eins oder Eins oder eine sigmoidale Funktion, welche auf dem Intervall $[0, 1]$ beschränkt ist. Die Aktivierungsfunktion die in der folgenden Analyse verwendet wird besitzt die Form

$$\phi(z) = \frac{1}{1 + \exp(-z)}. \quad (1)$$

Wie bereits erläutert wurde, sind biologische Neuronen Teil eines hierarchischen Netzwerkes, in dem das Signal von manchen in andere Neuronen eingespeist wird. Im Allgemeinen wird diese Struktur durch miteinander verbundene *Knoten* einer *Schicht* repräsentiert. Ein Netzwerk besteht aus einer Input-, ein oder mehrerer versteckter Schichten und einer Output-Schicht. Der Name der mittleren Schichten trägt der Tatsache Rechnung, dass das Wirken der versteckten Schichten zu einem gewissen Grad nur schwer zu beobachten, zuweilen sogar komplett unbeobachtbar ist. Für gewöhnlich sieht der Benutzer eines neuronalen Netzes nur das, was eingespeist wird, sowie dessen Ergebnis. Innerhalb eines *Knoten* wird die Aktivierungsfunktion des Neurons auf die gewichtete

Abbildung 1: Neuronales Netz mit einer versteckten Schicht.



Summe aller Kanten des Graphen angewandt. Der Output h_l^k eines *Knoten* der ersten versteckten Schicht $H_l^{(1)}$ ist somit als

$$\begin{aligned} h_l^{(k)} &= \phi(\beta_0 + x_1\beta_1 + \dots + x_{28}\beta_{28}) \\ &= \frac{1}{1 + \exp(-\beta_0 - \beta_1 w_1 - \dots - \beta_{28} x_{28})}. \end{aligned} \quad (2)$$

Das Signal o_1 des Output-Knotens O_1 ist wiederum eine gewichtete Summe der Outputs der versteckten Schicht.

$$o_1 = \phi \left(\sum_{q=0}^u \gamma_q h_q^{(1)} \right) \quad (3)$$

Mithilfe dieser Gleichungen lassen sich die Inputs vorwärts durch das Netzwerk propagieren, vorausgesetzt man kennt die Gewichte β und γ . Im Allgemeinen muss man die Gewichte zunächst ermitteln. Dies geschieht durch die Optimierung einer Kostenfunktion. Im Falle binärer Zustandsvariablen bietet sich die Kreuzentropie an.

$$J = -\frac{1}{n} \sum (y \log a + (1 - y) \log(1 - a)), \quad (4)$$

wobei a für die jeweilige Aktivierungsfunktion steht.

Die Kreuzentropie aus Gleichung 4 vereint zwei Eigenschaften, die sie als Kostenfunktion besonders auszeichnen. Sie ist größer Null und wenn der Output a des Endknotens nahe am tatsächlichen Zustand der Beobachtung liegt, ist der Wert des Summanden an dieser Stelle ebenfalls sehr gering. Für die optimalen Gewichte gilt es, J zu minimieren. Das Verfahren dazu nennt sich in der Literatur *backpropagation* und ist eine Anwendung der Kettenregel des Ableitens.

$$\begin{aligned}\frac{\partial J}{\partial \beta_j} &= -\frac{1}{n} \sum \left(\frac{y}{\sigma(z)} - \frac{1-y}{1-\sigma(z)} \right) \frac{\partial \sigma}{\partial \beta_j} \\ &= -\frac{1}{n} \sum \left(\frac{y}{\sigma(z)} - \frac{1-y}{1-\sigma(z)} \right) \sigma(z) (1-\sigma(z)) x_j.\end{aligned}\tag{5}$$

Wobei wir die sigmoidale Aktivierungsfunktion $\sigma(z)$ für a in Gleichung 4 eingesetzt haben. Gleichung 5 lässt sich weiter vereinfachen,

$$\begin{aligned}\frac{\partial J}{\partial \beta_j} &= \frac{1}{n} \sum \frac{\sigma(z) (1-\sigma(z)) x_j}{\sigma(z) (1-\sigma(z))} (\sigma(z) - y) \\ &= \frac{1}{n} \sum x_j (\sigma(z) - y) \stackrel{!}{=} 0.\end{aligned}\tag{6}$$

Der Ausdruck in 6 besagt, dass die Rate mit der das Gewicht gelernt wird, von dem Fehler in $(\sigma(z) - y)$ abhängt. Je größer dieser Fehler, desto schneller lernt das Neuron.

3.2 Architekturen und Hyperparameter

3.3 Ergebnisse

4 Zusammenfassung

Abbildung 2: Neuronales Netz mit zwei versteckten Schichten.

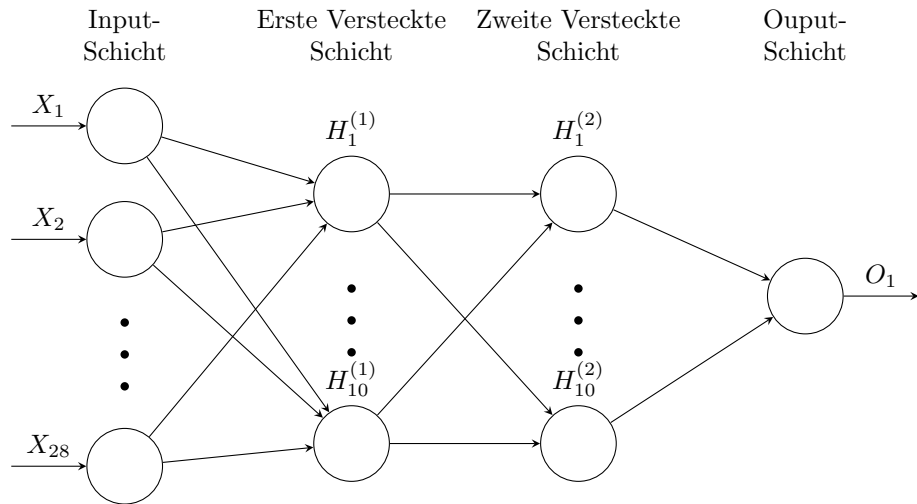


Abbildung 3: Neuronales Netz mit fünf versteckten Schichten.

