
Klassifikation der Schwierigkeitsgrade von Sudokus mit Methoden des maschinellen Lernens

Investigation of Sudoku difficulty levels
Bachelor-Thesis von Michael Bräunlein



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Fachbereich Informatik
Knowledge Engineering Group
Betreuer Prof. Dr. Johannes Fürnkranz

Zusammenfassung

Das Zahlenrätsel Sudoku ist weltweit bei Rätselliebhabern bekannt und beliebt. Seit seiner Veröffentlichung 1986 begeistern sich immer mehr Menschen für mehr oder weniger schwierige Exemplare. Sudokus finden sich im Internet, in der Rätsellecke der Tageszeitungen und sogar als ganze Bücher, um nur einige Erscheinungsorte zu nennen.

Die Regeln sind einfach zu lernen und doch kann man sich sehr lange mit Sudokus beschäftigen, da die schwersten Sudokus meist nur von Profis gelöst werden können.

Der Spielspass ist sehr stark davon abhängig, dass die Schwierigkeit zum persönlichen Können passt. Ist das Sudoku zu leicht, stellt es keine Herausforderung dar. Ist es zu schwer, kommt schnell ein Gefühl der Überforderung auf.

Die ausgewiesenen Schwierigkeitsstufen von Sudokus aus verschiedenen Quellen haben zwar oft die gleichen Namen wie zum Beispiel "Mittel", unterscheiden sich aber dennoch häufig nach Meinung des Spielers.

Das Ziel dieser Bachelorarbeit ist, Merkmale aus Sudokus zu extrahieren, anhand derer die Sudokus von einem Klassifizierer möglichst zuverlässig in Schwierigkeitsstufen eingeteilt werden können.

Inhaltsverzeichnis

1	Aufgabenstellung und Zielsetzung	3
2	Einführung	4
2.1	Die Regeln	4
2.2	Begriffserklärung	5
3	Lösungsmethoden	6
3.1	Kandidatenlisten	6
3.2	Full House	7
3.3	Naked Single	8
3.4	Hidden Single	9
3.5	Pointing Pair / Triple	10
3.6	Box-Line Reduction	11
3.7	Naked Subset	11
3.8	Hidden Subset	11
3.9	Fish	11
3.9.1	X-Wing	11
3.9.2	Swordfish	11
3.9.3	Jellyfish	11
3.10	Single Digit Patterns	11
3.10.1	Skyscarper	11
3.10.2	2-String Kite	11
3.10.3	Turbot Fish	11
3.10.4	Empty Rectangle	11
3.11	Wings	11
3.11.1	XY-Wing	11
3.11.2	XYZ-Wing	11
3.11.3	W-Wing	11
3.12	Sue de Coq	11
3.13	Coloring	11
3.14	Almost Locked Set	12
3.14.1	ALS XZ	12
3.14.2	ALS XY Wing	12
3.14.3	ALS Chain	12
4	Merkmalsextrahierung	13
4.1	Allgemeines Vorgehen	13
4.2	Entkopplung von konkreten Zahlen	13

1 Aufgabenstellung und Zielsetzung

Diese Bachelorarbeit beschäftigt sich mit der Einteilung von Sudokus in verschiedene Schwierigkeitsstufen. Hierzu sollen Methoden des maschinellen Lernens verwendet werden.

Es soll eine Methode gefunden werden, mit der Merkmale aus Sudokus extrahiert werden können, die dann als Feature Vektoren in einer .arff Datei¹ gesammelt werden. Die Feature Vektoren werden anschließend mit Hilfe von Weka² klassifiziert.

Es werden verschiedene Klassifikatoren und unterschiedliche Parameter betrachtet. Ausserdem werden Optimierungen der Featurevektoren diskutiert.

¹ <http://www.cs.waikato.ac.nz/ml/weka/arff.html>

² <http://www.cs.waikato.ac.nz/ml/weka/>

2 Einführung

Die Vorfahren des heutigen Sudokus waren vermutlich die lateinischen Quadrate, mit denen sich vor allem der Mathematiker Leonhard Euler befasste. Hier ging es darum, in ein Quadrat mit n Zeilen und n Spalten Symbole so einzutragen, dass jedes Symbol in jeder Spalte und Zeile jeweils genau einmal vorkommt.

1	2	3	4	5
2	5	4	1	3
3	4	5	2	1
4	3	1	5	2
5	1	2	3	4

Abbildung 2.1: Lateinisches Quadrat

Daraus hat sich das heutige Sudoku entwickelt, das sich nicht nur bei Mathematikern großer Beliebtheit erfreut.

2.1 Die Regeln

Diese Arbeit beschäftigt sich nur mit der meist verbreiteten Art von Sudoku. Dabei spielt man auf einem 9x9 Felder großen Spielfeld, das wiederum in neun 3x3 Felder große Blöcke eingeteilt ist. Weiter handelt es sich nur dann um ein Sudoku, wenn genau eine Lösung vorhanden ist. Ein Sudoku gilt dann als gelöst, wenn jede Zeile, jede Spalte und jeder Block die Ziffern 1 bis 9 genau einmal enthält.

4	8	1	5	6		3		
7	6	9	3	2	4		5	
3	5			7		6		
	9	7	2	8	5	2	1	
1			3		6	5	4	
5	4		1	3			8	6
	7	6	9	5	3		2	
2		5		4			3	
9	3	4					6	5

Abbildung 2.2: Sudoku

2.2 Begriffserklärung

Ein Sudoku besteht aus 81 *Feldern* oder *Zellen*. Diese bilden ein Quadrat der Größe 9x9, das *Grid*. Aufgrund dieser Aufteilung hat ein Sudoku 9 *Zeilen* und 9 *Spalten*. Das Grid wird in 9 Unterquadrate geteilt, die jeweils 3x3 Felder groß sind. Diese werden *Blöcke* genannt. Zeilen, Spalten und Blöcke werden unter dem Begriff *Figur* zusammengefasst. Die Nummerierung der Blöcke erfolgt zeilenweise von links oben nach rechts unten.

Vorgaben sind Zahlen, die schon von Anfang an gegeben sind.

In **Abbildung 2.2** sieht man im mittleren Block sogenannte *Kandidaten*. Ein Kandidat ist eine Zahl, die in der Zelle noch möglich ist. Jede Zelle hat ihre eigene Liste mit Kandidaten.

In der Beschreibung der Lösungstechniken ist es notwendig bestimmte Felder zu betrachten. Hierzu wird eine Abkürzung verwendet, die Zeile und Spalte enthält und somit eine Zelle eindeutig indentifiziert. z2s3 meint zum Beispiel die Zelle in Zeile 2 und Spalte 3.

In der folgenden Abbildung sind die erläuterten Begriffe zum besseren Verständniss eingetragen.

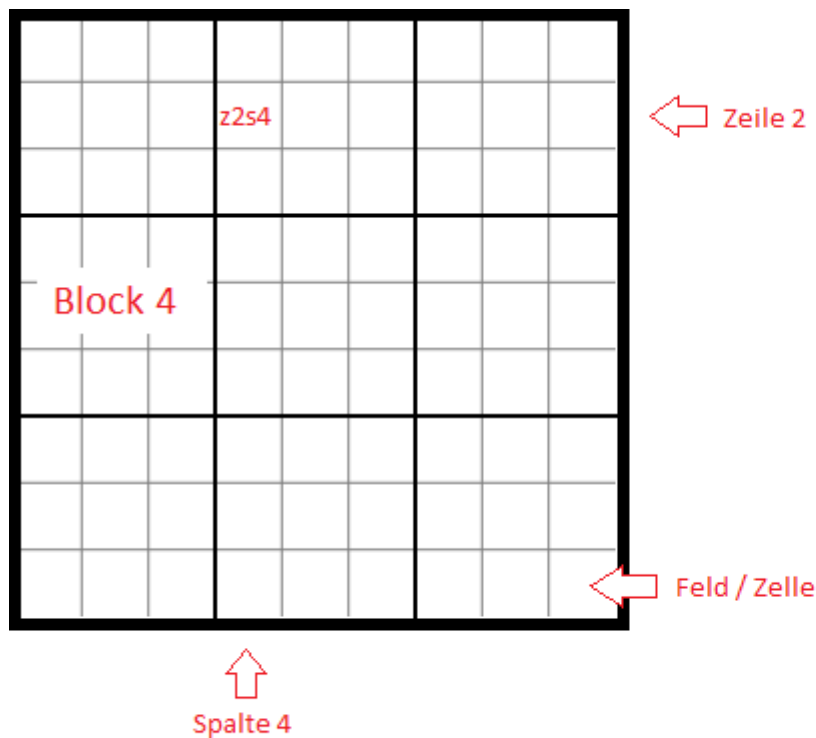


Abbildung 2.3: Begriffe

3 Lösungsmethoden

Alle in dieser Bachelorarbeit beschriebenen Techniken sind nicht im Rahmen dieser Arbeit entwickelt worden, sondern wurden aus verschiedenen Quellen zusammengetragen. Die Beschreibung der Lösungstechniken lehnt sich an die Beschreibung der Quellen an. Teile der Beispiele wurden aus den Quellen entnommen, dies ist entsprechend gekennzeichnet.

Grob kann man die Techniken zum Lösen von Sudokus in zwei Kategorien einteilen. Die erste Kategorie findet Zahlen heraus, die direkt in das Sudoku eingetragen werden können. Die Techniken der zweiten Kategorie entfernen Bedingungen in einzelnen Zellen des Sudokus.

3.1 Kandidatenlisten

Beim Lösen von Sudokus ist es üblich, in jedes Feld die Kandidaten einzutragen, die dort stehen können. Dabei wird vorerst nur die Sudoku Regel berücksichtigt, die besagt, dass in jeder Zeile die Zahlen 1 bis 9 vorkommen müssen. Wenn in einer Zeile nun die Zahl 3 vorkommt, dann kann sie in der selben Zeile nicht nochmal vorkommen, daher kann sie aus allen Kandidatenlisten der Zellen in der selben Zeile gelöscht werden. Dasselbe gilt für Spalten und Blöcke. Immer wenn eine Ziffer in ein Feld eingetragen wird, dann muss der Spieler die Liste der Kandidaten aktualisieren.

Kandidatenlisten sind keine eigene Lösungstechnik, sind aber wesentlicher Bestandteil vieler Techniken.

3.2 Full House

Wenn in einer Figur 8 Zahlen eingetragen sind, dann kann die Technik *Full House* angewendet werden. Da in jeder Figur die Zahlen 1 bis 9 stehen müssen, kann die fehlende Zahl einfach per Ausschluss ermittelt werden.

2	¹ 4 5 6	7	¹ 4 5 8	^{1 3} 5	³ 8	¹ 6	^{1 3} 4 5 6	^{1 3} 5 6 9
¹ 4 5	8	4 5 6	^{1 2} 4 5	9	^{2 3} 7	^{1 2} 6 7	^{1 3} 4 5 6 7	^{1 2 3} 5 6 7
¹ 4 5 9	3	4 5	6	¹ 5 7	² 7	8	¹ 4 5 7	^{1 2} 5 7 9
³ 5 7	⁵ 7	8	¹	6	4	9	¹ 7	^{1 2} 7
6	9	2	7	8	5	3	¹	4
⁴ 7	⁴ 7	1	3	2	⁹	5	^{7 8 6} 7 8 7	⁶
³ 4 5 7 8	^{4 5 6} 7	9	⁵ 8 7	^{5 3} 5	1	⁶ 7	2	^{5 6 3} 7 5 6
^{1 3} 5 7 8	^{1 2} 5 6 7	³ 5 6	² 5 8	4	^{2 3} 6 7 8	¹ 6 7	9	^{1 3} 5 6 7
^{1 3} 5 7	^{1 2} 5 6 7	³ 5 6	² 5 9 7	³ 5 7	^{2 3} 6 7 9	4	^{1 3} 5 6 7	8

Abbildung 3.1: Full House

In **Abbildung 3.1** fehlt in Zeile 5 nur noch eine Ziffer. Da die Zahlen 2 bis 9 bereits vorhanden sind, kann in das Feld z5s8 die Zahl 1 eingetragen werden.

3.3 Naked Single

Bei der Technik *Naked Single* werden Kandidatenlisten verwendet. Diese Technik kann angewendet werden, wenn in der Kandidatenliste einer Zelle nur noch ein Kandidat steht. Dieser Kandidat kann dann in die Zelle eingetragen werden. Das funktioniert aufgrund des Aufbaus der Kandidatenlisten. Diese enthalten zuerst alle Kandidaten und es werden immer dann Kandidaten entfernt, wenn dieser Kandidat nicht mehr als Ziffer in der Zelle stehen könnte weil er dort eine Regel verletzen würde. Wenn also nur noch ein Kandidat in der Kandidatenliste steht, dann bedeutet das, dass dieser Kandidat die einzige Ziffer zwischen 1 und 9 ist, die in der Zelle stehen kann ohne eine Regel zu verletzen.

<div>2 5 6 8</div>	1	<div>2 3 5 6</div>	9	<div>5 6</div>	<div>5 6</div>	7	4	<div>2 5 8</div>
<div>2 4 5 6 9</div>	<div>2 4 5 9</div>	<div>2 5 6</div>	8	<div>4 5 6</div>	<div>1 4 5 6 7</div>	<div>1 5</div>	<div>1 5</div>	3
<div>4 5 8</div>	7	5	3	2	<div>1 4 5</div>	6	9	<div>1 5 8</div>
<div>1 5 6 7 9</div>	<div>5 9</div>	4	<div>5 7</div>	3	<div>5 7 9</div>	2	<div>1 5 6 7 8</div>	<div>1 5 6 7 8 9</div>
<div>1 5 7 9</div>	<div>3 5 7 9</div>	<div>3 5</div>	6	<div>4 5 8 9</div>	2	<div>1 4 5 8 9</div>	<div>1 5 7 8</div>	<div>1 4 5 7 8 9</div>
<div>2 5 6 7 9</div>	<div>2 5 9</div>	8	<div>4 5 7</div>	1	<div>4 5 7 9</div>	3	<div>5 6 7</div>	<div>4 5 6 7 9</div>
<div>2 4 5</div>	8	1	<div>2 4 5</div>	7	<div>4 5 6 9</div>	<div>4 5 9</div>	3	<div>4 5 6 9</div>
3	<div>2 4 5</div>	<div>2 5 7</div>	<div>1 2 4 5</div>	<div>4 5 6 9</div>	8	<div>1 4 5 9</div>	<div>1 5 6 7</div>	<div>1 4 5 6 7 9</div>
<div>4 5 7</div>	6	9	<div>1 4 5</div>	<div>4 5</div>	3	<div>1 4 5 8</div>	2	<div>1 4 5 7 8</div>

Abbildung 3.2: Naked Single

Im oben stehenden Beispiel **Abbildung 3.2** sieht man sofort, dass die Kandidatenliste in z3s3 nur noch einen Eintrag enthält. Dieser kann nun einfach eingetragen werden.

3.4 Hidden Single

Auch die Technck *Hidden Single* arbeitet wieder mit Kandidatenlisten. Wenn in einer Figur eine Kandidatenliste die einzige ist, in der eine bestimmte Zahl vorkommt, dann kann diese Zahl direkt in die Zelle eingetragen werden. Wenn in dieser Zelle die Zahl nicht stünde, dann gäbe es in der Figur keine Möglichkeit mehr, dass die Zahl auftaucht und damit wäre die Sudoku Regel verletzt, nach der jede Zahl genau einmal enthalten sein muss.

4 5 3 9	2	8	4 5 6 1 9 5	7	4 6 3 9 9	6 4 9 3
4 5 3 9	1	6	4 5 8 9	3	4 2 9	7 4 9 2
4 3 7 9	4 7	4 3 9	4 6 9	2	4 8 9	5 1
1	3	7	2	9	5 8	4 5 6 8 6 4 5 6 8 8
4 5 6 8 9	4 5 6 8	4 5 2 9	7	3	1 5 8	4 5 6 8 9 1 2 8 9 2
5 8 9	5 8	2 5 9	1 5 5	4	6	3 1 2 8 9 7
2	9	1 3 4 5	1 4 5	7	1 4 5	5 6 3 8 6 5 6 8
7 5 3 7	5 5 7	5 3 5	8	6	2 5 9	1 4 2 3 5 9
4 5 6 8	4 5 6 8	4 5 1 5	3	1 5 5	1 2 4 5 9	7 2 6 8 9 5 6 8 9

Abbildung 3.3: Hidden Single

In **Abbildung 3.3** sieht man, dass die Zahl 6 in der Zeile 3 nur in z3s4 vorkommen kann. Daher kann man sie dort eintragen.

3.5 Pointing Pair / Triple

Bei der Technik *Pointing Pair / Triple* müssen zum ersten mal die Kandidatenlisten mehrerer Felder gleichzeitig betrachtet werden, was diese Technik etwas schwerer macht. Ausserdem ist diese Technik die erste, die Kandidaten aus Kandidatenlisten entfernt und nur bedingt zum Einsetzen von Zahlen in das Sudoku führt.

Es werden die Kandidatenlisten in Blöcken jeweils zeilen- und spaltenweise betrachtet. Die Technik *Pointing Pair / Triple* kann angewendet werden, wenn in einem Block eine Kandidat nur in Kandidatenlisten der selben Zeile oder Spalte vorkommt. Dann kann jedes weitere vorkommen der Zahl in einer Kandidatenliste der selben Zeile oder Spalte entfernt werden. Das gilt, da die Zahl genau einmal in dem Block vorkommen muss. Da alle möglichen Vorkommen der Zahl in der selben Zeile oder Spalte liegen ist klar, dass die Zahl in dieser Zeile oder Spalte vorkommt. Da sie aber kein zweites mal in der Zeile oder Spalte vorkommen darf muss sie aus den Kandidatenlisten entfernt werden, die nicht im selben Block liegen.

3	4	¹ 5 ₉	^{1 2} 5 ₈	⁵ 8 ₉	6	^{1 2} 5	7	^{1 2} 8
¹ 5 _{6 7}	8	¹ 5 ₆	^{1 2} 4 _{5 7}	^{4 5} 7	¹ 5 ₇	9	3	^{1 2} 4
¹ 5 _{7 9}	¹ 5 ₇	2	¹ 4 _{5 7 8}	3	¹ 5 _{7 9}	¹ 4 ₅	6	¹ 4 ₈
² 4 _{5 6 8}	^{2 3} 5 ₆	³ 4 _{5 6 8}	⁵ 7 ₈	1	⁵ 7 ₉	^{2 3} 4 ₆	² 4	^{2 3} 4 _{6 7 9}
^{1 2}	9	7	3	6	4	8	5	^{1 2}
¹ 4 _{5 6 8}	¹ 3 _{5 6}	¹ 3 _{4 5 6 8}	⁵ 7 ₈	⁵ 7 _{8 9}	2	¹ 3 _{4 6}	¹ 4	¹ 3 _{4 6 7 9}
^{1 2} 4 _{5 6 7 8 9}	^{1 2 3} 5 _{6 7 8 9}	¹ 3 _{4 5 6 8 9}	¹ 4 _{5 7}	^{4 5} 7	¹ 5 ₇	^{1 2 3} 4 ₆	^{1 2} 4	^{1 2 3} 4 ₆
^{1 2} 4 _{5 7}	^{1 2 3} 5 ₇	¹ 3 _{4 5}	6	^{4 5} 7	8	^{1 2 3} 4	9	^{1 2 3} 4
¹ 4 ₆	¹ 6 ₄	¹ 6 ₄	9	2	3	7	8	5

Abbildung 3.4: Pointing Pair / Triple

In **Abbildung 3.4** betrachten wir Block 8. Hier ist das Vorkommen der Zahl 1 in den Kandidatenlisten auf Zeile 7 beschränkt. Wie oben beschrieben können nun alle weiteren vorkommen in der selben Zeile, die nicht in Block 8 liegen aus den Kandidatenlisten entfernt werden. Im vorliegenden Beispiel führt das allerdings nicht dazu, dass eine neue Zahl in das Sudoku eingetragen wird. Dennoch ist das Sudoku nun genauer bestimmt, da weniger Möglichkeiten übrig sind.

3.6 Box-Line Reduction

3.7 Naked Subset

3.8 Hidden Subset

3.9 Fish

3.9.1 X-Wing

3.9.2 Swordfish

3.9.3 Jellyfish

3.10 Single Digit Patterns

3.10.1 Skyscarper

3.10.2 2-String Kite

3.10.3 Turbot Fish

3.10.4 Empty Rectangle

3.11 Wings

3.11.1 XY-Wing

3.11.2 XYZ-Wing

3.11.3 W-Wing

3.12 Sue de Coq

3.13 Coloring

3.14 Almost Locked Set

3.14.1 ALS XZ

3.14.2 ALS XY Wing

3.14.3 ALS Chain

4 Merkmalsextrahierung

4.1 Allgemeines Vorgehen

4.2 Entkopplung von konkreten Zahlen
