

Skript zu den Lehrveranstaltungen

THEORETISCHE INFORMATIK 1 + 2

Prof. Dr. Carsten Lutz

Basierend auf einem Skript von Franz Baader, TU Dresden

Inhaltsverzeichnis

Einführung	4
I. Endliche Automaten und reguläre Sprachen	9
1. Nichtdeterministische endliche Automaten	14
2. Deterministische endliche Automaten	20
3. Nachweis der Nichterkennbarkeit	30
4. Abschlusseigenschaften und Entscheidungsprobleme	33
5. Reguläre Ausdrücke und Sprachen	38
II. Grammatiken, kontextfreie Sprachen und Kellerautomaten	43
6. Die Chomsky-Hierarchie	44
7. Rechtslineare Grammatiken und reguläre Sprachen	47
8. Normalformen kontextfreier Grammatiken	50
9. Abschlusseigenschaften kontextfreier Sprachen	58
10. Kellerautomaten	62
III. Berechenbarkeit	69
11. Turingmaschinen	71
12. Zusammenhang zwischen Turingmaschinen und Grammatiken	80
13. Berechenbarkeit, Entscheidbarkeit, Aufzählbarkeit	87
14. Primitiv rekursive Funktionen und Loop-Programme	92
15. μ -rekursive Funktionen und While-Programme	104
16. Universelle Maschinen und unentscheidbare Probleme	111
17. Weitere unentscheidbare Probleme	119
IV. Komplexität	129
18. Komplexitätsklassen	130
19. NP-vollständige Probleme	133
Abkürzungsverzeichnis	142
Literatur	143

Hinweis

Dieses Skript ist als Hilfestellung für Studenten gedacht. Trotz großer Sorgfalt beim Erstellen kann keine Garantie für Fehlerfreiheit übernommen werden. Es wird explizit darauf hingewiesen, dass der prüfungsrelevante Stoff durch die Vorlesung bestimmt wird und mit dem Skriptinhalt nicht vollständig übereinstimmen muss.

Einführung

Die theoretische Informatik beschäftigt sich mit den mathematischen Grundlagen der Informatik und stellt ein wichtiges Fundament für zahlreiche andere Teilgebiete der Informatik dar. Insbesondere versucht man in der theoretischen Informatik, die zentralen Konzepte und Methoden der Informatik zu identifizieren und in abstrahierter Form zu studieren. Ein wesentlicher Schwerpunkt liegt dabei auf dem Erkennen und der exakten Beschreibung von prinzipiellen Grenzen der Informatik, wie z.B. den Grenzen der effizienten Berechenbarkeit.

Die theoretische Informatik ist in zahlreiche Teilgebiete untergliedert, wie etwa die Komplexitätstheorie, die Algorithmentheorie, die Kryptographie und die Datenbanktheorie. Die Lehrveranstaltungen „Theoretische Informatik 1 + 2“ geben eine Einführung in folgende zwei zentrale Bereiche der theoretischen Informatik:

Automatentheorie und formale Sprachen

Behandelt in Theoretische Informatik 1 / Teile I + II dieses Skriptes

Im Mittelpunkt stehen *Wörter* und *formale Sprachen* (Mengen von Wörtern). Diese sind ein nützliches Abstraktionsmittel in der Informatik. Man kann z.B. die Eingabe oder Ausgabe eines Programmes als Wort betrachten und die Menge der syntaktisch korrekten Eingaben als Sprache. Wichtige Fragestellungen sind z.B.:

- Was sind geeignete Beschreibungsmittel für (meist unendliche) formale Sprachen ? (z.B. mittels Automaten und Grammatiken)
- Was für verschiedene Typen von Sprachen lassen sich unterscheiden?
- Was für Eigenschaften haben die verschiedenen Sprachtypen?

Berechenbarkeit und Komplexität

Behandelt in Theoretische Informatik 2 / Teile III + IV dieses Skriptes

Hier geht es darum, welche Problem und Funktionen prinzipiell berechenbar sind und welche nicht. Ausserdem wird untersucht, welcher zeitliche Aufwand zur Berechnung eines Problems / einer Funktion notwendig ist (unabhängig vom konkreten Algorithmus). Wichtige Fragestellungen sind z.B.:

- Was kennzeichnet Berechenbarkeit (berechenbar *durch was?*)
- Gibt es Funktionen, die prinzipiell nicht berechenbar sind?
- Kann man jede berechenbare Funktion mit akzeptablem Zeit- und Speicherplatzaufwand berechnen?
- Für in der Informatik häufig auftretende Probleme/Funktionen: wie viel Zeit und Speicherplatz braucht man mindestens, also bei optimalem Algorithmus?

Automatentheorie und formale Sprachen

Formale Sprachen, also (endliche oder unendliche) Mengen von Wörtern, sind ein wichtiger Abstraktionsmechanismus der Informatik. Hier ein paar Anwendungsbeispiele:

- Die Menge aller wohlgeformten Programme in einer gegebenen Programmiersprache wie Pascal, Java, oder C++
- Die Menge aller wohlgeformten Eingaben für ein Programm oder eine Form auf einer Webseite (z.B. Menge aller Kontonummern / Menge aller Geburtsdaten)
- Jeder Suchausdruck (z.B. Linux Regular Expression) definiert eine formale Sprache
- Kommunikationsprotokolle: z.B. Menge aller wohlgeformten TCP-Pakete, wenn man die konkreten Nutzdaten wegabstrahiert
- Das “erlaubte” Verhalten von Soft- und Hardwaresystemen kann in sehr natürlicher Weise als formale Sprache modelliert werden.

Wir beginnen mit einem kurzen Überblick über die zentralen Betrachtungsgegenstände und Fragestellungen.

1. Charakterisierung:

Wie beschreibt man die (meist unendlichen) Mengen von Wörtern mit endlichem Aufwand?

- *Automaten* oder *Maschinen*, die genau die Elemente der Menge akzeptieren. Wir werden viele verschiedene Automatenmodelle kennenlernen, wie z.B. endliche Automaten, Kellerautomaten und Turingmaschinen.
- *Grammatiken*, welche genau die Elemente der Menge generieren; auch hier gibt es viele verschiedene Typen, z.B. rechtslineare Grammatiken und kontextfreie Grammatiken (vgl. auch VL „Praktische Informatik“: kontextfreie Grammatiken (EBNF) zur Beschreibung der Syntax von Programmiersprachen).
- *Ausdrücke*, welche beschreiben, wie man die Sprache aus Basissprachen mit Hilfe gewisser Operationen (z.B. Vereinigung) erzeugen kann.

Abhängig von dem verwendeten Automaten-/Grammatiktyp erhält man verschiedene Klassen von Sprachen. Wir werden hier die vier wichtigsten Klassen betrachten, welche in der **Chomsky-Hierarchie** zusammengefasst sind:

Klasse	Automatentyp	Grammatiktyp
Typ 0	Turingmaschine (TM)	allgemeine Chomsky-Grammatik
Typ 1	TM mit linearer Bandbeschränkung	kontextsensitive Grammatik
Typ 2	Kellerautomat	kontextfreie Grammatik
Typ 3	endlicher Automat	einseitig lineare Grammatik

Bei Typ 3 existiert auch eine Beschreibung durch reguläre Ausdrücke. Am wichtigsten sind die Typen 2 und 3; beispielsweise kann Typ 2 weitgehend die Syntax von Programmiersprachen beschreiben.

2. Welche **Problemstellungen** sind für eine Sprachklasse entscheidbar und mit welchem Aufwand? Die folgenden Probleme sind "klassisch":

- *Wortproblem*: gegeben eine Beschreibung der Sprache L (z.B. durch Automat, Grammatik, Ausdruck, ...) und ein Wort w . Gehört w zu L ?

Anwendungsbeispiele:

- Programmiersprache, deren Syntax durch eine kontextfreie Grammatik beschrieben ist. Entscheide für ein gegebenes Programm P , ob dieses syntaktisch korrekt ist.
- Suchpattern für Textdateien sind häufig reguläre Ausdrücke. Suche die Dateien (Wörter), welche das Suchpattern enthalten (zu der von ihm beschriebenen Sprache gehören).

- *Leerheitsproblem*: gegeben eine Beschreibung der Sprache L . Ist L leer?

Anwendungsbeispiel:

Wenn ein Suchpattern die leere Sprache beschreibt, so muss man die Dateien nicht durchsuchen, sondern kann ohne weiteren Aufwand melden, dass das Pattern nicht sinnvoll ist.

- *Äquivalenzproblem*: Beschreiben zwei verschiedene Beschreibungen dieselbe Sprache?

Anwendungsbeispiel:

Jemand vereinfacht die Grammatik einer Programmiersprache, um sie übersichtlicher zu gestalten. Beschreibt die vereinfachte Grammatik wirklich dieselbe Sprache wie die ursprüngliche?

3. Welche **Abschlusseigenschaften** hat eine Sprachklasse?

z.B. Abschluss unter Durchschnitt, Vereinigung und Komplement: wenn L_1, L_2 in der Sprachklasse enthalten, sind es dann auch $L_1 \cap L_2, L_1 \cup L_2, \overline{L_1}$?

Anwendungsbeispiele:

- Suchpattern: Suche nach Dateien, die das Pattern *nicht* enthalten (Komplement) oder die zwei Pattern enthalten (Durchschnitt).
- Reduziere das Äquivalenzproblem auf das Leerheitsproblem, ohne die gewählte Klasse von Sprachen zu verlassen: Statt „ $L_1 = L_2$?“ entscheidet man, ob $(L_1 \cap \overline{L_2}) \cup (L_2 \cap \overline{L_1})$ leer ist.

I. Endliche Automaten und reguläre Sprachen

Einführung

Es werden einige der grundlegenden Begriffe der Vorlesung “Theoretische Informatik 1” eingeführt sowie ein erster, informeller Blick auf endliche Automaten geworfen.

Formale Sprachen

Die zentralen Begriffe der Vorlesung “Theoretische Informatik 1” sind die eines Wortes und einer formalen Sprache.

Alphabet. Ein Alphabet ist eine endliche Menge von Symbolen. Beispiele sind:

- $\Sigma_1 = \{a, b, c, \dots, z\}$;
- $\Sigma_2 = \{0, 1\}$;
- $\Sigma_3 = \{0, \dots, 9\} \cup \{, \}$;
- $\Sigma_4 = \{ \text{program, const, var, label, procedure, function, type, begin, end, if, then, else, case, of, repeat, until, while, do, for, to} \} \cup \{ \text{VAR, VALUE, FUNCTION} \}$

Als Platzhalter für Alphabetssymbole benutzen wir in der Regel a, b, c, \dots . Alphabete bezeichnen wir meist mit Σ .

Obwohl die Symbole von Σ_4 aus mehreren Buchstaben der üblichen Schriftsprache bestehen, betrachten wir sie doch als unteilbare Symbole. Die Elemente von Σ_4 sind genau die Schlüsselworte der Programmiersprache Pascal. Konkrete Variablennamen, Werte und Funktionsaufrufe sind zu den Schlüsselworten VAR, VALUE, FUNCTION abstrahiert um Endlichkeit des Alphabetes zu gewährleisten.

Wort. Ein Wort ist eine endliche Folge von Symbolen. Ein Wort $w = a_1 \dots a_n$ mit $a_i \in \Sigma$ heißt Wort *über dem Alphabet* Σ . Beispiele sind:

- $w = abc$ ist ein Wort über Σ_1 ;
- $w = 1000110$ ist ein Wort über Σ_2 ;
- $w = 10,0221,4292,,$ ist ein Wort über Σ_3 ;

- Jedes Pascalprogramm kann als Wort über Σ_4 betrachtet werden, wenn man jede konkrete Variable durch das Schlüsselwort VAR ersetzt, jeden Wert durch VALUE und jeden Funktionsaufruf durch FUNCTION.

Als Platzhalter für Worte verwenden wir meist w, v, u . Die Länge eines Wortes w wird mit $|w|$ bezeichnet, es gilt also z.B. $|aba| = 3$. Manchmal ist es praktisch, auch die Anzahl Vorkommen eines Symbols a in einem Wort w in kurzer Weise beschreiben zu können. Wir verwenden hierfür $|w|_a$, es gilt also z.B. $|aba|_a = 2$, $|aba|_b = 1$, $|aba|_c = 0$. Einen Spezialfall stellt das *leere Wort* dar, also die leere Folge von Symbolen. Dieses wird durch ε bezeichnet. Es ist das einzige Wort mit $|w| = 0$.

Formale Sprache. Eine (formale) Sprache ist eine Menge von Wörtern. Mit Σ^* bezeichnen wir die Sprache, die aus *allen* Wörtern über dem Alphabet Σ bestehen. Eine Sprache $L \subseteq \Sigma^*$ heißt *Sprache über dem Alphabet Σ* . Beispiele sind:

- $L = \emptyset$
- $L = \{abc\}$
- $L = \{a, b, c, ab, ac, bc\}$
- $L = \{w \in \{a, \dots, z\}^* \mid w \text{ ist ein Wort der deutschen Sprache} \}$
- L als Menge aller Worte über Σ_4 , die wohlgeformte Pascal-Programme beschreiben

Als Platzhalter für Sprachen verwenden wir meist L . Beachten Sie, dass Sprachen sowohl endlich als auch unendlich sein können. Interessant sind meist nur unendliche Sprachen. Als nützliche Abkürzung führen wir Σ^+ für die Menge $\Sigma^* \setminus \{\varepsilon\}$ aller nicht-leeren Wörter über Σ ein.

Operationen auf Sprachen und Wörtern

Im folgenden werden wir sehr viel mit Wörtern und formalen Sprachen umgehen. Dazu verwenden wir in erster Linie die folgenden Operationen.

Präfix, Suffix, Infix: Zu den natürlichsten und einfachsten Operationen auf Wörtern gehört das Bilden von Präfixen, Suffixen und Infixen:

$$\begin{array}{lll} u \text{ ist Präfix von } v & \text{gdw.} & v = uw \text{ für ein } w \in \Sigma^*. \\ u \text{ ist Suffix von } v & \text{gdw.} & v = wu \text{ für ein } w \in \Sigma^*. \\ u \text{ ist Infix von } v & \text{gdw.} & v = w_1uw_2 \text{ für } w_1, w_2 \in \Sigma^*. \end{array}$$

Die Präfixe von $aabbcc$ sind also beispielsweise $a, aa, aab, aabb, aabb, aabbcc$. Dieses Wort hat 21 Infixe.

Konkatenation: Eine Operation, die auf Wörtern sowie auf Sprachen angewendet werden kann. Auf Wörtern u und v bezeichnet die Konkatenation $u \cdot v$ das Wort uv , das man durch einfaches "Hintereinanderschreiben" erhält. Es gilt also z.B.

$abb \cdot ab = abbab$. Auf Sprachen bezeichnet die Konkatenation das Hintereinanderschreiben *beliebiger* Worte aus den beteiligten Sprachen:

$$L_1 \cdot L_2 := \{u \cdot v \mid (u \in L_1) \wedge (v \in L_2)\}$$

Es gilt also z.B.

$$\{aa, a\} \cdot \{ab, b, aba\} = \{aaab, aab, aaaba, ab, aaba\}.$$

Sowohl auf Sprachen als auch auf Wörtern wird der Konkatenationspunkt häufig weggelassen, wir schreiben also z.B. L_1L_2 statt $L_1 \cdot L_2$.

Man beachte, dass $\emptyset \cdot L = L \cdot \emptyset = \emptyset$. Konkatenation ist assoziativ, es gilt also $(L_1 \cdot L_2) \cdot L_3 = L_1 \cdot (L_2 \cdot L_3)$. Sie ist nicht kommutativ, im allgemeinen gilt also nicht $L_1 \cdot L_2 = L_2 \cdot L_1$.

Boolesche Operationen: Es handelt sich um die üblichen Booleschen Mengenoperationen, angewendet auf formale Sprachen:

$$\begin{array}{ll} \text{Vereinigung} & L_1 \cup L_2 := \{w \mid w \in L_1 \text{ oder } w \in L_2\} \\ \text{Durchschnitt} & L_1 \cap L_2 := \{w \mid w \in L_1 \text{ und } w \in L_2\} \\ \text{Komplement} & \overline{L_1} := \{w \mid w \in \Sigma^* \wedge w \notin L_1\} \end{array}$$

Manchmal verwenden wir zusätzlich die Differenz, also

$$L_1 \setminus L_2 := L_1 \cap \overline{L_2} = \{w \mid w \in L_1 \wedge w \notin L_2\}.$$

Vereinigung und Durchschnitt sind sowohl assoziativ als auch kommutativ.

Kleene-Stern: Der Kleene-Stern bezeichnet die beliebig (aber nur endlich) oft iterierte Konkatenation. Gegeben eine Sprache L definiert man zunächst induktive Sprachen L^0, L^1, \dots und darauf basierend dann die Anwendung des Kleene-Sterns erhaltene Sprache L^* :

$$\begin{array}{ll} L^0 & := \{\varepsilon\} \\ L^{n+1} & := L^n \cdot L \\ L^* & := \bigcup_{n \geq 0} L^n \end{array}$$

Für $L = \{a, ab\}$ gilt also z.B. $L^0 = \{\varepsilon\}$, $L^1 = L$, $L^2 = \{aa, aab, aba, abab\}$, etc. Offensichtlich ist L^* unendlich gdw. $L \neq \emptyset$.

Man beachte, dass das leere Wort per Definition *immer* in L^* enthalten ist, unabhängig davon, was L für eine Sprache ist. Manchmal verwenden wir auch die Variante ohne das leere Wort:

$$L^+ := \bigcup_{n \geq 1} L^n = L^* \setminus \{\varepsilon\}.$$

Einige einfache Beobachtungen sind $\emptyset^* = \{\varepsilon\}$, $(L^*)^* = L^*$ und $L^* \cdot L^* = L^*$. Es ist hier wichtig, \emptyset (die leere Sprache), $\{\varepsilon\}$ (die Sprache, die das leere Wort enthält) und ε (das leere Wort) sorgsam auseinander zu halten.

Weitere nützliche Notation:

- $|w|$ bezeichnet die Länge des Wortes w (=Anzahl Symbole in w),
z.B. $|aabba| = 5$
- a^n bezeichnet das Wort, das aus n mal dem Symbol a besteht,
z.B. $a^3 = aaa, b^5 = bbbbb$
die formale Definition ist induktiv: $a^0 = \varepsilon$ und $a^{i+1} = a^i \cdot a$
- die gerade eingeführte Notation wird auch für Wörter w verwendet:
z.B. $(abc)^3 = abcabcabc$ (aber $abc^3 = abccc$)
- $|w|_a$ bezeichnet die Anzahl Vorkommen des Symbols a im Wort w ,
z.B. $|abbaa|_a = 3$ und $|abbaa|_b = 2$

Endliche Automaten

Endliche Automaten stellen ein einfaches und dennoch sehr nützliches Mittel zum Beschreiben von formalen Sprachen dar. Die charakteristischen Merkmale eines endlichen Automaten sind

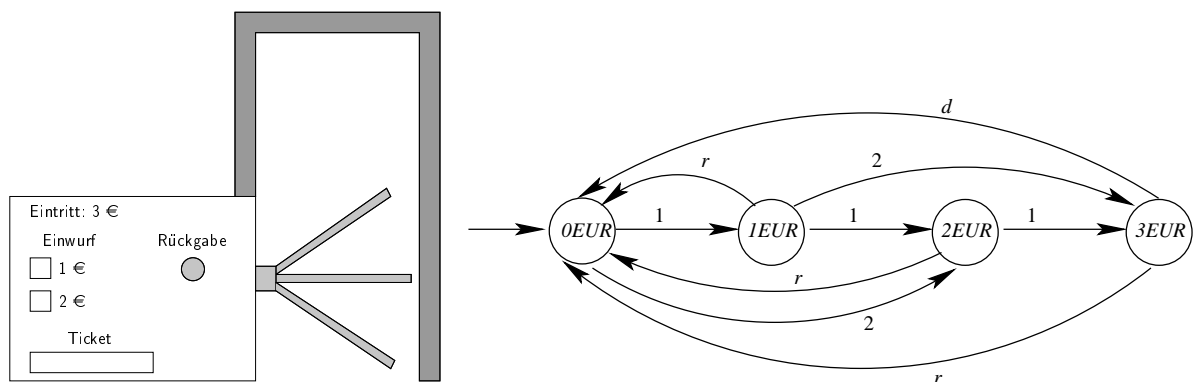
- eine endliche Menge von (internen) Zuständen, in denen sich der Automat befinden kann
- in einer festen Tabelle beschriebene Übergänge zwischen Zuständen; der Folgezustand ist dabei abhängig vom momentanen Zustand und der Eingabe

Ein endlicher Automat kann sowohl einen mechanischen Automaten als auch einen elektronischen Automaten darstellen. Ursprünglich wurden endliche Automaten sogar als einfaches Modell des menschlichen Gehirns verwendet. Als erstes Beispiel verwenden wir einen Eintrittsautomat mit Drehsperre.

Beispiel: (Eintrittsautomat)

Eingabe: 1, 2, r, d (*r*: Geldrückgabe; *d*: Drehsperre dreht sich)

Zustände: 0EUR, 1EUR, 2EUR, 3EUR



Der dargestellte Automat regelt eine Drehsperre. Es können Münzen im Wert von 1 oder 2 Euro eingeworfen werden. Nach Einwurf von 3 Euro wird die insgesamt Arretierung der Drehsperre gelöst und der Eintritt freigegeben. Der Automat gibt kein Wechselgeld zurück sondern nimmt einen zu hohen Betrag einfach nicht an (Münzen fallen durch). Man kann jederzeit den Rückgabeknopf drücken, um den bereits gezahlten Betrag zurückzuerhalten.

In der schematischen Darstellung kennzeichnen die Kreise die internen Zustände und die Pfeile die Übergänge. Die Pfeilbeschriftung gibt die jeweilige Eingabe an, unter der der Übergang erfolgt. Man beachte, dass

- nur der Zustand 3EUR einen Übergang vom Typ d erlaubt. Dadurch wird modelliert, dass nur durch Einwurf von 3,- Euro der Eintritt ermöglicht wird.
- das Drehen der Sperre als Eingabe angesehen wird. Man könnte dies auch als Ausgabe modellieren. Wir werden in dieser Vorlesung jedoch keine endlichen Automaten mit Ausgabe (sogenannte Transduktoren) betrachten.

Die Übergänge können als festes Programm betrachtet werden, das der Automat ausführt.

Man beachte den offensichtlichen Zusammenhang zu formalen Sprachen: die Menge der möglichen Eingaben bildet ein Alphabet. Jede (Gesamt-)Eingabe des Automaten ist ein Wort über dem Alphabet. Wenn man 3EUR als Zielzustand betrachtet, so bildet die Menge der Eingaben, mittels derer dieser Zustand erreicht werden kann, eine (unendliche) formale Sprache.

Randbemerkung.

Im Prinzip sind Rechner ebenfalls endliche Automaten: Sie haben nur endlich viel Speicherplatz und daher nur eine endliche Menge möglicher Konfigurationen (Prozessorzustand + Belegung der Speicherzellen). Die Konfigurationsübergänge werden bestimmt durch Verdrahtung und Eingaben (Tastatur, Peripheriegeräte).

Wegen der extrem großen Anzahl von Zuständen sind endliche Automaten aber keine geeignete Abstraktion für Rechner. Ausserdem verwendet man einen Rechner (z.B. bei der Programmierung) nicht als endlichen Automat indem man z.B. ausnutzt, dass der Arbeitsspeicher ganz genau 2GB gross ist. Stattdessen nimmt man den Speicher als potentiell unendlich an und verlässt sich auf Techniken wie Swapping und Paging. In einer geeigneten Abstraktion von Rechnern sollte daher auch der Speicher als unendlich angenommen werden. Das wichtigste solche Modell ist die Turingmaschine, die wir später im Detail kennenlernen werden.

1. Nichtdeterministische endliche Automaten

Wir werden verschiedene Begriffe von endlichen Automaten kennenlernen und später sehen, dass sie alle dieselbe Sprachklasse definieren. Zunächst führen wir jedoch den sehr allgemeinen Begriff des *Transitionssystems* ein. Es handelt sich hier noch nicht um einen endlichen Automaten, da ein Transitionssystem auch unendlich sein darf.

Definition 1.1 (Transitionssystem)

Ein *Transitionssystem* ist von der Form $\mathcal{A} = (Q, \Sigma, I, \Delta, F)$, wobei

- Q eine Menge von Zuständen ist (nicht notwendigerweise endlich!),
- Σ ein Alphabet ist,
- $I \subseteq Q$ eine Menge von Anfangszuständen ist,
- $\Delta \subseteq Q \times \Sigma \times Q$ eine Übergangsrelation (Transitionsrelation) ist,
- $F \subseteq Q$ eine Menge von Endzuständen ist.

Ein *endlicher Automat* unterscheidet sich von einem Transitionssystem dadurch, dass er nur endlich viele Zustände und nur einen Anfangszustand hat.

Definition 1.2 (nichtdeterministischer endlicher Automat)

Ein Transitionssystem $\mathcal{A} = (Q, \Sigma, I, \Delta, F)$ nennt man *nichtdeterministischer endlicher Automat (NEA)*, wenn

- $|Q| < \infty$
- $|I| = 1$

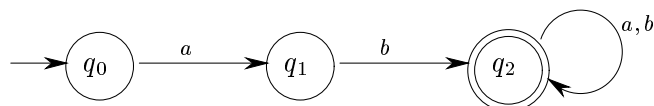
Anstelle von $(Q, \Sigma, \{q_0\}, \Delta, F)$ schreiben wir $(Q, \Sigma, q_0, \Delta, F)$.

Beispiel 1.3

Der NEA $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ mit den Komponenten

- $Q = \{q_0, q_1, q_2\}$
- $\Sigma = \{a, b\}$
- $F = \{q_2\}$
- $\Delta = \{(q_0, a, q_1), (q_1, b, q_2), (q_2, a, q_2), (q_2, b, q_2)\}$

wird graphisch dargestellt als:



Wie im obigen Beispiel werden wir Automaten häufig als kantenbeschriftete Graphen darstellen, wobei die Zustände des Automaten die Knoten des Graphen sind und die Übergänge als Kanten gesehen werden (beschriftet mit einem Alphabetsymbol). Der Startzustand wird durch einen Pfeil gekennzeichnet und die Endzustände durch einen Doppelkreis.

Ein Transitionssystem $\mathcal{A} = (Q, \Sigma, I, \Delta, F)$ enthält Wörter über dem Alphabet Σ als Eingabe. Ein Übergang $(q, a, q') \in \Delta$ bedeutet, dass der Automat im Zustand q und bei Eingabe a in den Zustand q' wechseln kann. Intuitiv liest der Automat also das Eingabewort von links nach rechts. Er akzeptiert die Eingabe, wenn dabei der Endzustand erreicht werden kann und verwirft sie sonst. Dies werden wir im folgenden formal definieren.

Definition 1.4 (Pfad)

Ein *Pfad* in einem Transitionssystem \mathcal{A} ist eine Folge

$$\pi = (p_0, a_1, p_1)(p_1, a_2, p_2) \cdots (p_{n-1}, a_n, p_n)$$

mit $(p_i, a_{i+1}, p_{i+1}) \in \Delta$ für $i = 0, \dots, n - 1$.

π ist ein Pfad von p nach q , wenn $p = p_0$ und $q = p_n$.

Die *Beschriftung* des Pfades π ist das Wort $\beta(\pi) := a_1 \cdots a_n$.

Die *Länge* des Pfades π ist n .

Für $n = 0$ sprechen wir vom *leeren Pfad*, welcher die Beschriftung ε hat.

Die Existenz eines Pfades in \mathcal{A} von p nach q mit Beschriftung w beschreiben wir durch

$$p \xrightarrow{\mathcal{A}}_w q.$$

Für Mengen $Q_1, Q_2 \subseteq Q$ heißt $Q_1 \xrightarrow{\mathcal{A}}_w Q_2$, dass es Zustände $p_1 \in Q_1$ und $p_2 \in Q_2$ gibt mit $p_1 \xrightarrow{\mathcal{A}}_w p_2$.

Für den NEA aus Beispiel 1.3 gilt beispielsweise für alle $w \in \{a, b\}^*$: $q_0 \xrightarrow{\mathcal{A}}_{abw} q_2$.

Definition 1.5 (Akzeptiertes Wort, erkannte Sprache)

Das Transitionssystem $\mathcal{A} = (Q, \Sigma, I, \Delta, F)$ *akzeptiert* das Wort $w \in \Sigma^*$ gdw. $I \xrightarrow{\mathcal{A}}_w F$ gilt. Die von \mathcal{A} *erkannte Sprache* ist $L(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ akzeptiert } w\}$.

Für den NEA aus Beispiel 1.3 ist $L(\mathcal{A}) = \{abw \mid w \in \Sigma^*\}$, d.h. \mathcal{A} akzeptiert genau die Wörter über $\Sigma = \{a, b\}$, welche mit ab beginnen.

Definition 1.6 (Erkennbarkeit einer Sprache)

Eine Sprache $L \subseteq \Sigma^*$ heißt *erkennbar*, wenn es einen NEA \mathcal{A} gibt mit $L = L(\mathcal{A})$.

Da wir uns bei einem Transitionssystem i.a. nur für die erkannten Sprachen interessieren, bezeichnen wir zwei Transitionssysteme als *äquivalent*, wenn sie dieselbe Sprache akzeptieren.

Manchmal ist es bequem, ein Transitionssystem oder einen NEA als Graph aufzufassen. Wir werden diesen Zusammenhang darum kurz etwas näher beleuchten.

Definition (Gerichteter Graph, Erreichbarkeit)

Ein *gerichteter Graph* $G = (V, E)$ besteht aus einer Menge V von *Knoten* und einer Menge $E \subseteq V \times V$ von *Kanten*. Ein Knoten $v \in V$ ist *erreichbar* von einem Knoten $u \in V$ wenn es eine Knotenfolge v_1, \dots, v_n gibt, mit $n \geq 1$, so dass $v_1 = u$, $v_n = v$, und $(v_i, v_{i+1}) \in E$ für $1 \leq i < n$.

Graphen sind eine fundamentale Struktur sowohl in der Informatik als auch in der Mathematik. Sie können z.B. zur Repräsentation von Kommunikationsnetzen und Datenstrukturen verwendet werden. Ein Transitionssystem $\mathcal{A} = (Q, \Sigma, I, \Delta, F)$ kann als gerichteter Graph (V, E) dargestellt werden mit

- $V = Q$ und
- $E = \{(q, q') \mid \exists a \in \Sigma : (q, a, q') \in \Delta\}$.

Man fasst also die Zustände als Knoten und die Übergänge als Kanten auf, wobei die Alphabetsymbole sowie Start- und Endzustände wegabstrahiert werden. Es ist natürlich auch möglich, nur solche Übergänge, die bestimmte, ausgewählte Symbole betreffen, als Kanten in den Graph aufzunehmen.

Das *Erreichbarkeitsproblem* für gerichtete Graphen ist das folgende Problem: gegeben ein endlicher gerichteter Graph $G = (V, E)$ und zwei Knoten $u_0, v_0 \in V$, entscheide ob v_0 von u_0 erreichbar ist. Ein einfacher Algorithmus, um dieses Problem zu lösen, ist der folgende: gegeben $G = (V, E)$ und u_0, v_0 , berechne eine Folge von Knotenmengen $V_0 \subseteq V_1 \subseteq V_2 \dots$ so dass $V_0 = \{u_0\}$ und $V_i \subseteq V$ für alle $i \geq 0$ durch Setzen von

$$V_{i+1} := V_i \cup \{v' \mid \exists v \in V_i : (v, v') \in E\}$$

Die Folge wird Schritt für Schritt berechnet. Sobald $V_i = V_{i+1}$ gilt, stoppt der Algorithmus, gibt „ja“ zurück wenn $v_0 \in V_i$ und „nein“ sonst. Wir werden in der Übung zeigen, dass der Algorithmus tatsächlich Erreichbarkeit entscheidet.

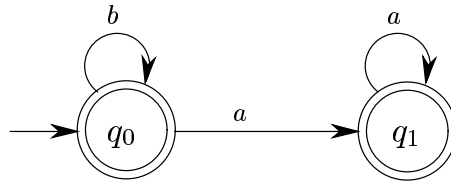
Theorem (Komplexität und Entscheidbarkeit von Erreichbarkeit)

Das Erreichbarkeitsproblem für gerichtete Graphen ist effektiv in Linearzeit entscheidbar, das heisst mit $\mathcal{O}(|V| + |E|)$ Schritten.

Der oben vorgestellte, einfache Algorithmus läuft nicht in Linearzeit sondern in quadratischer Zeit, macht also $\mathcal{O}((|V| + |E|)^2)$ Schritte. Die Vorstellung eines Linearzeitalgorithmus für das Erreichbarkeitsproblem ist in dieser Vorlesung nicht möglich und findet sich z.B. im Buch „Introduction to Algorithms“ von Cormen, Leiserson, Rivest und Stein.

Beispiel 1.7

$L = \{w \in \{a, b\}^* \mid ab \text{ ist nicht Infix von } w\}$ ist erkennbar, denn L wird z.B. von dem folgenden Automaten erkannt:



Für manche Konstruktionen ist es günstig, auch endliche Automaten zu betrachten, bei denen Übergänge mit Wörtern beschriftet sind.

Definition 1.8 (NEA mit Wortübergängen, ε -NEA)

Ein NEA mit Wortübergängen hat die Form $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$, wobei Q, Σ, q_0, F wie beim NEA definiert sind und $\Delta \subseteq Q \times \Sigma^* \times Q$ eine endliche Menge von Wortübergängen ist.

Ein ε -NEA ist ein NEA mit Wortübergängen, wobei für alle $(q, w, q') \in \Delta$ gilt: $|w| \leq 1$ (d.h. $w \in \Sigma$ oder $w = \varepsilon$).

Pfade, Pfadbeschriftungen und erkannte Sprache werden entsprechend wie für NEAs definiert. Zum Beispiel hat der Pfad $(q_0, ab, q_1)(q_1, \varepsilon, q_2)(q_2, bb, q_3)$ die Beschriftung $ab \cdot \varepsilon \cdot bb = abbb$.

Satz 1.9

Zu jedem NEA mit Wortübergängen kann man effektiv einen äquivalenten NEA konstruieren.

Effektiv bedeutet hier, dass es einen Algorithmus gibt, der als Eingabe einen NEA mit Wortübergängen erhält und als Ausgabe einen äquivalenten NEA ohne Wortübergänge liefert.

Man zeigt Satz 1.9 mit Umweg über ε -NEAs.

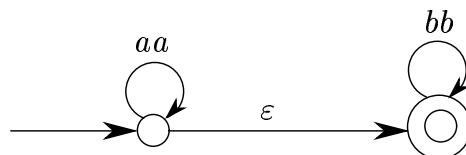
Lemma 1.10

Zu jedem NEA mit Wortübergängen kann man effektiv einen äquivalenten ε -NEA konstruieren.

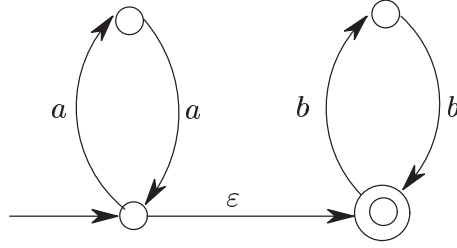
Beweis. Man ersetzt jeden Wortübergang $(q, a_1 \cdots a_n, q')$ mit $n > 1$ durch Symbolübergänge $(q, a_1, p_1), (p_1, a_2, p_2), \dots, (p_{n-1}, a_n, q')$, wobei p_1, \dots, p_{n-1} jeweils neue Hilfszustände sind (die nicht zur Endzustandsmenge dazugenommen werden). Man sieht leicht, dass dies einen äquivalenten ε -NEA liefert. \square

Beispiel 1.11

Der NEA mit Wortübergängen, der durch die folgende Darstellung gegeben ist:



wird überführt in einen äquivalenten ε -NEA:



Lemma 1.12

Zu jedem ε -NEA kann man effektiv einen äquivalenten NEA konstruieren.

Beweis. Der ε -NEA $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ sei gegeben. Wir konstruieren daraus einen NEA \mathcal{A}' ohne ε -Übergänge wie folgt:

$\mathcal{A}' = (Q, \Sigma, q_0, \Delta', F')$, wobei

- $\Delta' := \left\{ (p, a, q) \in Q \times \Sigma \times Q \mid p \xrightarrow{a}_{\mathcal{A}} q \right\}$
- $F' := \begin{cases} F \cup \{q_0\} & \text{falls } q_0 \xrightarrow{\varepsilon}_{\mathcal{A}} F \\ F & \text{sonst} \end{cases}$

Noch zu zeigen: $L(\mathcal{A}) = L(\mathcal{A}')$ und \mathcal{A}' kann effektiv bestimmt werden.

1. $L(\mathcal{A}') \subseteq L(\mathcal{A})$:

Sei $w = a_1 \cdots a_n \in L(\mathcal{A}')$. Dann gibt es in \mathcal{A}' Pfad

$$(p_0, a_1, p_1)(p_1, a_2, p_2) \cdots (p_{n-1}, a_n, p_n) \quad \text{mit } p_0 = q_0, p_n \in F'$$

Nach Definition von Δ' gibt es auch in \mathcal{A} einen Pfad π von p_0 nach p_n mit Beschriftung w .

1. Fall: $p_n \in F$

Dann zeigt π , dass $w \in L(\mathcal{A})$.

2. Fall: $p_n \in F' \setminus F$, d.h. $p_n = q_0$

Nach Definition von F' gilt $p_n = q_0 \xrightarrow{\varepsilon}_{\mathcal{A}} \hat{p}$ für ein $\hat{p} \in F$. Es gibt also in \mathcal{A} einen Pfad von p_0 nach \hat{p} mit Beschriftung w , daher $w \in L(\mathcal{A})$.

2. $L(\mathcal{A}) \subseteq L(\mathcal{A}')$:

Sei $w \in L(\mathcal{A})$ und $\pi = (p_0, a_1, p_1)(p_1, a_2, p_2) \cdots (p_{m-1}, a_m, p_m)$ ein Pfad in \mathcal{A} mit $p_0 = q_0, p_m \in F$ und Beschriftung w , wobei $a_i \in \Sigma \cup \{\varepsilon\}$. Seien i_1, \dots, i_n die Indizes mit $a_{i_j} \neq \varepsilon$. Offenbar ist $w = a_{i_1} \cdots a_{i_n}$.

1. Fall: $n > 0$, also $w \neq \varepsilon$

Nach Definition von Δ' ist

$$(p_0, a_{i_1}, p_{i_1})(p_{i_1}, a_{i_2}, p_{i_2}) \cdots (p_{i_{n-1}}, a_{i_n}, p_m)$$

ein nicht-leerer Pfad in \mathcal{A}' . Aus $p_m \in F$ folgt $p_m \in F'$, was $w \in L(\mathcal{A}')$ zeigt.

2. Fall: $n = 0$, also $w = \varepsilon$

Dann ist $a_1 = \dots = a_m = \varepsilon$. Es gilt also $q_0 = p_0 \xrightarrow{\varepsilon}_{\mathcal{A}} p_m \in F$, was $q_0 \in F'$ liefert. Also $w = \varepsilon \in L(\mathcal{A}')$.

3. Δ' und F' können effektiv bestimmt werden:

- $p \xrightarrow{a}_{\mathcal{A}} q$ gilt genau dann, wenn es Zustände $p', q' \in Q$ gibt, für die gilt:

$$p \xrightarrow{\varepsilon}_{\mathcal{A}} p', \quad (p', a, q') \in \Delta, \quad q' \xrightarrow{\varepsilon}_{\mathcal{A}} q$$

Man muss nur endlich viele p', q' prüfen. Ob „ $(p', a, q') \in \Delta$?“ kann effektiv geprüft werden, da Δ endliche Menge. Weiterhin sind „ $p \xrightarrow{\varepsilon}_{\mathcal{A}} p'$?“ sowie „ $q' \xrightarrow{\varepsilon}_{\mathcal{A}} q$?“ Erreichbarkeitsprobleme in dem endlichen Graphen $G = (V, E)$ mit $V = Q$ und

$$E = \{(u, v) \mid (u, \varepsilon, v) \in \Delta\},$$

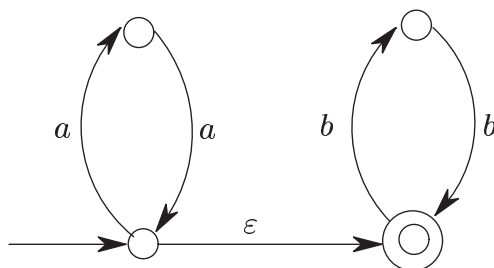
können also effektiv entschieden werden.

- „ $q_0 \xrightarrow{\varepsilon}_{\mathcal{A}} F$?“ ist ebenfalls Erreichbarkeitsproblem.

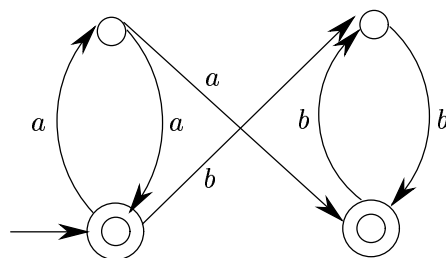
□

Beispiel: (zu Lemma 1.12)

Der ε -NEA aus Beispiel 1.11



wird in folgenden NEA überführt:



Lemma 1.12 kann man auch verwenden, wenn man zeigen will, dass es zu jedem *endlichen* Transitionssystem $\mathcal{A} = (Q, \Sigma, I, \Delta, F)$ einen äquivalenten NEA \mathcal{A}' gibt. Man definiert für einen neuen Zustand $q_0 \notin Q$ einen ε -NEA $\mathcal{A}'' = (Q \cup \{q_0\}, \Sigma, q_0, \Delta'', F)$ mit

$$\Delta'' = \Delta \cup \{(q_0, \varepsilon, q) \mid q \in I\}$$

\mathcal{A}'' ist äquivalent zu \mathcal{A} . Man erhält nun \mathcal{A}' aus \mathcal{A}'' wie im Beweis des Lemmas beschrieben.

2. Deterministische endliche Automaten

Die bisher betrachteten NEAs heißen nichtdeterministisch, weil es in ihnen zu einem Paar $(q, a) \in Q \times \Sigma$ mehr als einen Übergang geben kann, d.h. $q', q'' \in Q$ mit $q' \neq q''$ und $(q, a, q') \in \Delta$, $(q, a, q'') \in \Delta$.

Definition 2.1 (deterministischer endlicher Automat)

Ein NEA $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ heißt *deterministisch* (DEA), falls es für alle $q \in Q$, $a \in \Sigma$ genau ein $q' \in Q$ gibt mit $(q, a, q') \in \Delta$.

Anstelle der Übergangsrelation Δ verwenden wir dann die *Übergangsfunktion*

$$\delta : Q \times \Sigma \rightarrow Q$$

mit $\delta(q, a) = q'$ gdw. $(q, a, q') \in \Delta$. DEAs werden in der Form $(Q, \Sigma, q_0, \delta, F)$ geschrieben.

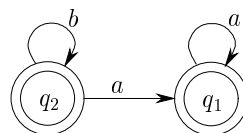
Beachte:

Zusätzlich zur *Eindeutigkeit* eines Übergangs haben wir auch die *Existenz* eines Übergangs für jedes Paar (q, a) gefordert.

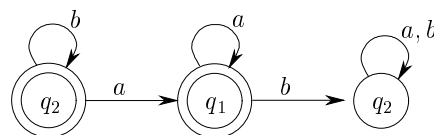
In der Literatur wird für deterministische Automaten manchmal nur die Eindeutigkeit des Übergangs gefordert, d.h. für (q, a) darf es *höchstens ein* q' mit $(q, a, q') \in \Delta$ geben. Man spricht dann bei Existenz aller Übergänge von *vollständigen* Automaten, d.h. wenn es für jedes (q, a) *mindestens ein* q' gibt mit $(q, a, q') \in \Delta$. Bezüglich der erkannten Sprachklasse macht dies keinen Unterschied.

Beispiel 2.2

Der NEA



erfüllt zwar die Eindeutigkeitsanforderung, er ist aber kein DEA in unserem Sinn, da für (q_1, b) kein Übergang existiert. Einen äquivalenten DEA erhält man durch Hinzunahme eines „Papierkorbzustandes“:



Definition 2.3 (kanonische Fortsetzung von δ)

Die *kanonische Fortsetzung* von $\delta : Q \times \Sigma \rightarrow Q$ auf eine Funktion $\delta^* : Q \times \Sigma^* \rightarrow Q$ wird induktiv (über die Wortlänge) definiert:

- $\delta^*(q, \varepsilon) := q$
- $\delta^*(q, wa) := \delta(\delta^*(q, w), a)$

Der Einfachheit halber werden wir in Zukunft auch für Wörter w die Schreibweise $\delta(q, w)$ statt $\delta^*(q, w)$ verwenden.

Beachte:

Für einen DEA $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ gilt:

- a) $\delta(q, w) = q'$ gdw. q' ist der eindeutige Zustand mit $q \xrightarrow{w}_{\mathcal{A}} q'$.
- b) $L(\mathcal{A}) = \{w \in \Sigma^* \mid \delta(q_0, w) \in F\}$
- c) $\delta(q, uv) = \delta(\delta(q, u), v)$

Wir zeigen nun, dass die bei der Definition von DEAs gemachten Einschränkungen (Existenz und Eindeutigkeit von Übergängen) keine echten Einschränkungen sind:

Satz 2.4 (Rabin/Scott)

Zu jedem NEA kann man effektiv einen äquivalenten DEA konstruieren.

Damit folgt, dass NEAs und DEAs dieselbe Klasse von Sprachen akzeptieren. Bevor wir den Beweis dieses Satzes angeben, skizzieren wir kurz die

Beweisidee:

Der Beweis dieses Satzes verwendet die sogenannte *Potenzmengenkonstruktion*: die Zustandsmenge des DEA ist die Potenzmenge 2^Q der Zustandsmenge Q des NEA.

Sei $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ ein NEA. Um für $w \in \Sigma^*$ zu entscheiden, ob $w \in L(\mathcal{A})$ ist, betrachtet man $\{q \in Q \mid q_0 \xrightarrow{w}_{\mathcal{A}} q\} \in 2^Q$. Es ist $w \in L(\mathcal{A})$ gdw. diese Menge enthält mindestens einen Endzustand.

Wir definieren einen DEA mit Zustandsmenge 2^Q und Übergangsfunktion δ so, dass $\delta(\{q_0\}, w) = \{q \mid q_0 \xrightarrow{w}_{\mathcal{A}} q\}$.

Beweis. Sei der NEA $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ gegeben. Der DEA $\mathcal{A}' = (2^Q, \Sigma, \{q_0\}, \delta, F')$ ist definiert durch:

- $\delta(P, a) = \bigcup_{p \in P} \{p' \mid (p, a, p') \in \Delta\}$ für alle $P \in 2^Q$ und $a \in \Sigma$
- $F' = \{P \in 2^Q \mid P \cap F \neq \emptyset\}$

Wir benötigen im Folgenden die

Hilfsaussage: $q' \in \delta(\{q\}, w)$ gdw. $q \xrightarrow{w}_{\mathcal{A}} q'$ (★)

Daraus folgt $L(\mathcal{A}) = L(\mathcal{A}')$, da:

$$\begin{aligned}
 w \in L(\mathcal{A}) & \quad \text{gdw.} \quad \exists q \in F : q_0 \xrightarrow{w}_{\mathcal{A}} q && (\text{Def. } L(\mathcal{A})) \\
 & \quad \text{gdw.} \quad \exists q \in F : q \in \delta(\{q_0\}, w) && (\text{Hilfsaussage}) \\
 & \quad \text{gdw.} \quad \delta(\{q_0\}, w) \cap F \neq \emptyset \\
 & \quad \text{gdw.} \quad \delta(\{q_0\}, w) \in F' && (\text{Def. } F') \\
 & \quad \text{gdw.} \quad w \in L(\mathcal{A}')
 \end{aligned}$$

Beweis der Hilfsaussage mittels Induktion über $|w|$:

Induktionsanfang: $|w| = 0$

$$q' \in \delta(\{q\}, \varepsilon) \quad \text{gdw.} \quad q = q' \quad \text{gdw.} \quad q \xrightarrow{\varepsilon}_{\mathcal{A}} q'$$

Induktionsannahme: Die Hilfsaussage ist bereits gezeigt für alle $w \in \Sigma^*$ mit $|w| \leq n$

Induktionsschritt: $|w| = n + 1$

Sei $w = ua$ mit $u \in \Sigma^*$, $|u| = n$ und $a \in \Sigma$. Es gilt:

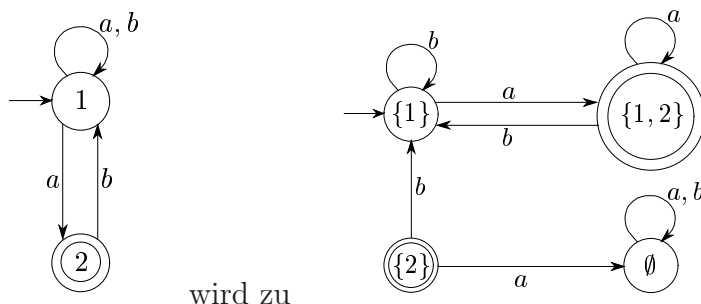
$$\begin{aligned}
 \delta(\{q\}, ua) &= \delta(\delta(\{q\}, u), a) && (\text{Def. 2.3}) \\
 &= \bigcup_{q' \in \delta(\{q\}, u)} \{q'' \mid (q', a, q'') \in \Delta\} && (\text{Def. } \delta) \\
 &= \bigcup_{q \xrightarrow{u}_{\mathcal{A}} q'} \{q'' \mid (q', a, q'') \in \Delta\} && (\text{Ind.Voraus.}) \\
 &= \{q'' \mid q \xrightarrow{ua}_{\mathcal{A}} q''\} && (\text{Def. Pfad})
 \end{aligned}$$

Daraus folgt sofort die Hilfsaussage für $w = ua$.

□

Beispiel 2.5

Der NEA \mathcal{A} (links) wird mit der Potenzmengenkonstruktion transformiert in den DEA \mathcal{A}' (rechts):



Nachteilig an dieser Konstruktion ist, dass die Zustandsmenge *exponentiell* vergrößert wird. Im Allgemeinen kann man dies nicht vermeiden, in manchen Fällen kommt man aber doch mit weniger Zuständen aus. Wir werden im Folgenden ein Verfahren angeben, welches zu einem gegebenen DEA einen äquivalenten DEA mit minimaler Zustandszahl, also einen reduzierten DEA konstruiert. Dieses Verfahren besteht aus 2 Schritten:

1. Schritt: Eliminieren unerreichbarer Zustände

Definition 2.6 (Erreichbarkeit eines Zustandes)

Ein Zustand q des DEA $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ heißt *erreichbar*, falls es ein Wort $w \in \Sigma^*$ gibt mit $\delta(q_0, w) = q$. Sonst heißt q *unerreichbar*.

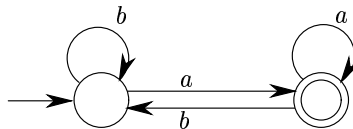
Da für die erkannte Sprache nur Zustände wichtig sind, welche von q_0 erreicht werden, erhält man durch Weglassen unerreichbarer Zustände einen äquivalenten Automaten:

$\mathcal{A}_0 = (Q_0, \Sigma, q_0, \delta_0, F_0)$ mit

- $Q_0 = \{q \in Q \mid q \text{ ist erreichbar}\}$
- $\delta_0 = \delta \upharpoonright_{Q_0 \times \Sigma}$ Beachte: Für $q \in Q_0$ und $a \in \Sigma$ ist auch $\delta(q, a) \in Q_0$!
- $F_0 = F \cap Q_0$

Beispiel 2.5 (Fortsetzung)

Im Automaten \mathcal{A}' von Beispiel 2.5 sind die Zustände $\{2\}$ und \emptyset nicht erreichbar. Durch Weglassen dieser Zustände erhält man den DEA \mathcal{A}'_0 :



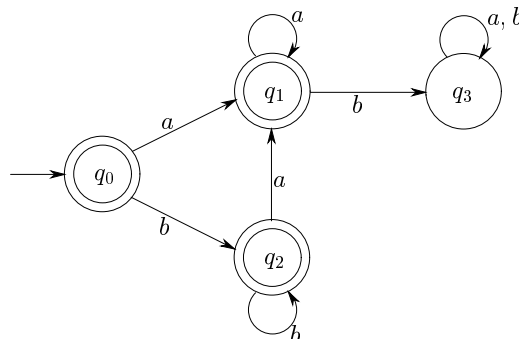
Bei der Potenzmengenkonstruktion kann man die Betrachtung unerreichbarer Elemente von 2^Q vermeiden, indem man mit $\{q_0\}$ beginnt und nur die davon erreichbaren Mengen sukzessive konstruiert.

2. Schritt: Zusammenfassen äquivalenter Zustände

Ein DEA ohne unerreichbare Zustände muss noch nicht minimal sein, da er noch verschiedene Zustände enthalten kann, die sich „gleich“ verhalten in Bezug auf die erkannte Sprache.

Beispiel 2.7

Der folgende DEA erkennt dieselbe Sprache wie der DEA aus Beispiel 2.2, hat aber einen Zustand mehr. Dies kommt daher, dass q_0 und q_2 äquivalent sind.



Definition 2.8 (Äquivalenz von Zuständen)

Es sei $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ ein DEA. Für $q \in Q$ sei $\mathcal{A}_q = (Q, \Sigma, q, \delta, F)$. Zwei Zustände $q, q' \in Q$ heißen \mathcal{A} -äquivalent ($q \sim_{\mathcal{A}} q'$) gdw. $L(\mathcal{A}_q) = L(\mathcal{A}_{q'})$.

Lemma 2.9

- 1) $\sim_{\mathcal{A}}$ ist eine Äquivalenzrelation auf Q , d.h. reflexiv, transitiv und symmetrisch.
- 2) $\sim_{\mathcal{A}}$ ist verträglich mit der Übergangsfunktion, d.h.

$$q \sim_{\mathcal{A}} q' \Rightarrow \forall a \in \Sigma : \delta(q, a) \sim_{\mathcal{A}} \delta(q', a)$$

- 3) $\sim_{\mathcal{A}}$ kann effektiv bestimmt werden.

Beweis.

- 1) ist klar, da „ $=$ “ reflexiv, transitiv und symmetrisch ist.
- 2) lässt sich wie folgt herleiten:

$$\begin{aligned} q \sim_{\mathcal{A}} q' &\Rightarrow L(\mathcal{A}_q) = L(\mathcal{A}_{q'}) \\ &\Rightarrow \forall w \in \Sigma^* : \delta(q, w) \in F \Leftrightarrow \delta(q', w) \in F \\ &\Rightarrow \forall a \in \Sigma \forall v \in \Sigma^* : \delta(q, av) \in F \Leftrightarrow \delta(q', av) \in F \\ &\Rightarrow \forall a \in \Sigma \forall v \in \Sigma^* : \delta(\delta(q, a), v) \in F \Leftrightarrow \delta(\delta(q', a), v) \in F \\ &\Rightarrow \forall a \in \Sigma : L(\mathcal{A}_{\delta(q, a)}) = L(\mathcal{A}_{\delta(q', a)}) \\ &\Rightarrow \forall a \in \Sigma : \delta(q, a) \sim_{\mathcal{A}} \delta(q', a) \end{aligned}$$

- 3) Wir definieren Approximationen \sim_k von $\sim_{\mathcal{A}}$:

- $q \sim_0 q'$ gdw. $q \in F \Leftrightarrow q' \in F$
- $q \sim_{k+1} q'$ gdw. $q \sim_k q'$ und $\forall a \in \Sigma : \delta(q, a) \sim_k \delta(q', a)$

Es gilt $Q \times Q \supseteq \sim_0 \supseteq \sim_1 \supseteq \sim_2 \supseteq \dots \supseteq \sim_{\mathcal{A}}$. Da Q endlich ist, gibt es ein k mit $\sim_k = \sim_{k+1}$. Man zeigt leicht, dass dann $\sim_k = \sim_{\mathcal{A}}$ ist.

□

Die $\sim_{\mathcal{A}}$ -Äquivalenzklasse eines Zustands $q \in Q$ ist $[q]_{\mathcal{A}} := \{q' \in Q \mid q \sim_{\mathcal{A}} q'\}$.

Beispiel 2.7 (Fortsetzung)

Für den Automaten aus Beispiel 2.7 gilt:

- \sim_0 hat die Klassen $F = \{q_0, q_1, q_2\}$ und $Q \setminus F = \{q_3\}$.
- \sim_1 hat die Klassen $\{q_1\}, \{q_0, q_2\}, \{q_3\}$.
Zum Beispiel ist $\delta(q_0, b) = \delta(q_2, b) \in F$ und $\delta(q_1, b) \notin F$.
- $\sim_2 = \sim_1 = \sim_{\mathcal{A}}$.

Definition 2.10 (Quotientenautomat)

Der *Quotientenautomat* $\tilde{\mathcal{A}} = (\tilde{Q}, \Sigma, [q_0]_{\mathcal{A}}, \tilde{\delta}, \tilde{F})$ zu $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ ist definiert durch:

- $\tilde{Q} := \{[q]_{\mathcal{A}} \mid q \in Q\}$
- $\tilde{\delta}([q]_{\mathcal{A}}, a) := [\delta(q, a)]_{\mathcal{A}}$ (repräsentantenunabhängig wegen Lemma 2.9)
- $\tilde{F} := \{[q]_{\mathcal{A}} \mid q \in F\}$

Lemma 2.11

$\tilde{\mathcal{A}}$ ist äquivalent zu \mathcal{A} .

Beweis. Es ist einfach, per Induktion über $|w|$ zu zeigen, dass die Erweiterung von $\tilde{\delta}$ auf Wörter sich wie folgt verhält:

$$\tilde{\delta}([q_0], w) \in \tilde{F} \quad \underline{\text{gdw.}} \quad [\delta(q_0, w)]_{\mathcal{A}} \in \tilde{F} \quad (*)$$

Nun gilt:

$$\begin{aligned} w \in L(\mathcal{A}) & \quad \underline{\text{gdw.}} \quad \delta(q_0, w) \in F \\ & \quad \underline{\text{gdw.}} \quad [\delta(q_0, w)]_{\mathcal{A}} \in \tilde{F} \quad (\text{Def. } \tilde{F}) \\ & \quad \underline{\text{gdw.}} \quad \tilde{\delta}([q_0]_{\mathcal{A}}, w) \in \tilde{F} \quad (*) \\ & \quad \underline{\text{gdw.}} \quad w \in L(\tilde{\mathcal{A}}) \end{aligned}$$

□

Definition 2.12 (reduzierter Automat zu einem DEA)

Für einen DEA \mathcal{A} bezeichnet $\mathcal{A}_{red} := \tilde{\mathcal{A}}_0$ den *reduzierten Automaten*, den man aus \mathcal{A} durch Eliminieren unerreichbarer Zustände und Zusammenfassen äquivalenter Zustände erhält.

Wir werden im Folgenden zeigen:

- \mathcal{A}_{red} kann nicht weiter vereinfacht werden, d.h. er ist der kleinste DEA, mit dem man $L(\mathcal{A})$ akzeptieren kann.
- \mathcal{A}_{red} hängt nur von $L(\mathcal{A})$ und nicht von \mathcal{A} ab, d.h. gilt $L(\mathcal{A}) = L(\mathcal{B})$, so auch $\mathcal{A}_{red} \simeq \mathcal{B}_{red}$ (gleich bis auf Zustandsumbenennung, isomorph).

Dazu konstruieren wir zu einer erkennbaren Sprache L einen „kanonischen“ Automaten \mathcal{A}_L und zeigen, dass jeder reduzierte Automat, der L erkennt, isomorph zu \mathcal{A}_L ist.

Definition 2.13 (Nerode-Rechtskongruenz)

Es sei $L \subseteq \Sigma^*$ eine beliebige Sprache. Für $u, v \in \Sigma^*$ definieren wir:
 $u \simeq_L v \quad \underline{\text{gdw.}} \quad \forall w \in \Sigma^* : uw \in L \Leftrightarrow vw \in L$.

Beispiel 2.14

Wir betrachten die Sprache

$$L = \{w \in \{a, b\}^* \mid ab \text{ ist nicht Infix von } w\}$$

(vgl. Beispiele 1.7, 2.2, 2.7)

- Es gilt:

$\varepsilon \simeq_L b :$	$\forall w : \varepsilon w \in L$	<u>gdw.</u>	$w \in L$
		<u>gdw.</u>	w enthält ab nicht
		<u>gdw.</u>	bw enthält ab nicht
		<u>gdw.</u>	$bw \in L$
- $\varepsilon \not\simeq_L a :$ $\varepsilon b \in L$, aber $a \cdot b \notin L$

Lemma 2.15 (Eigenschaften von \simeq_L)

- 1) \simeq_L ist eine Äquivalenzrelation.
- 2) \simeq_L ist Rechtskongruenz, d.h. zusätzlich zu 1) gilt: $u \simeq_L v \Rightarrow \forall w \in \Sigma^* : uw \simeq_L vw$.
- 3) L ist Vereinigung von \simeq_L -Klassen:

$$L = \bigcup_{u \in L} [u]_L$$

wobei $[u]_L := \{v \mid u \simeq_L v\}$.

- 4) Ist $L = L(\mathcal{A})$ für einen DEA \mathcal{A} , so ist die Anzahl der \simeq_L -Klassen \leq der Zustandszahl von \mathcal{A} . (Die Anzahl der \simeq_L -Klassen heißt Index von \simeq_L .)

Beweis.

- 1) folgt aus der Definition von \simeq_L , da „ \Leftrightarrow “ reflexiv, transitiv und symmetrisch ist.
- 2) Damit $uw \simeq_L vw$ gilt, muss für alle $w' \in \Sigma^*$ gelten:

(\star) $uww' \in L \Leftrightarrow vww' \in L$

 Nun ist aber $ww' \in \Sigma^*$ und daher folgt (\star) aus $u \simeq_L v$.
- 3) Zeige $L = \bigcup_{u \in L} [u]_L$.

„ \subseteq “: klar, da $u \in [u]_L$.

„ \supseteq “: Sei $u \in L$ und $v \in [u]_L$.

Wegen $\varepsilon \in \Sigma^*$ folgt aus $u = u \cdot \varepsilon \in L$ und $v \simeq_L u$ auch $v = v \cdot \varepsilon \in L$.

- 4) Es sei $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ ein DEA mit $L = L(\mathcal{A})$.

Hilfsaussage: Aus $\delta(q_0, u) = \delta(q_0, v)$ folgt $u \simeq_L v$, denn:

$\forall w : uw \in L$	<u>gdw.</u>	$\delta(q_0, uw) \in F$
	<u>gdw.</u>	$\delta(\delta(q_0, u), w) \in F$
	<u>gdw.</u>	$\delta(\delta(q_0, v), w) \in F$
	<u>gdw.</u>	$\delta(q_0, vw) \in F$
	<u>gdw.</u>	$vw \in L$

Also gibt es maximal so viele verschiedene Klassen, wie es verschiedene Zustände gibt (Schubfachprinzip).

□

Beispiel 2.14 (Fortsetzung)

\simeq_L hat drei Klassen:

- $[\varepsilon]_L = \{b\}^*$
- $[a]_L = \{b\}^* \cdot \{a\}^+$
- $[ab]_L = \Sigma^* \cdot \{ab\} \cdot \Sigma^*$

Man kann die \simeq_L -Klassen nun als Zustände eines Automaten für L verwenden.

Definition 2.16 (Transitionssystem \mathcal{A}_L zu einer Sprache L)

Zu $L \subseteq \Sigma^*$ ist das Transitionssystem $\mathcal{A}_L := (Q', \Sigma, q'_0, \delta', F')$ definiert durch:

- $Q' := \{[u]_L \mid u \in \Sigma^*\}$
- $q'_0 := [\varepsilon]_L$
- $\delta'([u]_L, a) := [ua]_L$ (repräsentantenunabhängig wegen Lemma 2.15, 2))
- $F' := \{[u]_L \mid u \in L\}$

Lemma 2.17

Hat \simeq_L endlichen Index, so ist \mathcal{A}_L ein DEA mit $L = L(\mathcal{A}_L)$.

Beweis. Hat \simeq_L endlich viele Klassen, so hat \mathcal{A}_L endlich viele Zustände. Außerdem gilt:

$$\begin{aligned} L(\mathcal{A}_L) &= \{u \mid \delta'(q'_0, u) \in F'\} \\ &= \{u \mid \delta'([\varepsilon]_L, u) \in F'\} \\ &= \{u \mid [u]_L \in F'\} \quad (\text{wegen } \delta'([u]_L, v) = [uv]_L) \\ &= \{u \mid u \in L\} \\ &= L \end{aligned}$$

□

Satz 2.18 (Satz von Nerode)

Eine Sprache L ist erkennbar gdw. \simeq_L hat endlichen Index (d.h. endlich viele Klassen).

Beweis.

„ \Rightarrow “: Ergibt sich unmittelbar aus Lemma 2.15, 4).

„ \Leftarrow “: Ergibt sich unmittelbar aus Lemma 2.17, da \mathcal{A}_L DEA ist, der L erkennt.

□

Beispiel 2.19 (nichterkennbare Sprache)

Die Sprache $L = \{a^n b^n \mid n \geq 0\}$ ist nicht erkennbar, da für $n \neq m$ gilt: $a^n \not\sim_L a^m$. In der Tat gilt $a^n b^n \in L$, aber $a^m b^n \notin L$. Daher hat \simeq_L unendlichen Index.

Wir werden im nächsten Abschnitt noch eine weitere Möglichkeit sehen, wie man für eine Sprache zeigen kann, dass sie *nicht* regulär ist.

Zunächst untersuchen wir aber den Zusammenhang zwischen *reduzierten Automaten* und der *Nerode-Rechtskongruenz*.

Definition 2.20 (isomorph)

Zwei DEAs $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ und $\mathcal{A}' = (Q', \Sigma, q'_0, \delta', F')$ sind *isomorph* ($\mathcal{A} \simeq \mathcal{A}'$) gdw. es eine Bijektion $\pi : Q \rightarrow Q'$ gibt mit:

- $\pi(q_0) = q'_0$
- $\pi(F) = F'$, wobei $\pi(F) := \{\pi(q) \mid q \in F\}$
- $\pi(\delta(q, a)) = \delta'(\pi(q), a)$ für alle $q \in Q, a \in \Sigma$

Lemma 2.21

$\mathcal{A} \simeq \mathcal{A}' \Rightarrow L(\mathcal{A}) = L(\mathcal{A}')$

Beweis. Es sei $\pi : Q \rightarrow Q'$ der Isomorphismus. Durch Induktion über $|w|$ zeigt man leicht, dass $\pi(\delta(q, w)) = \delta'(\pi(q), w)$. Daher gilt:

$$\begin{aligned} w \in L(\mathcal{A}) \quad \text{gdw.} \quad & \delta(q_0, w) \in F \\ & \text{gdw.} \quad \pi(\delta(q_0, w)) \in F' \quad (\text{wegen } F' = \pi(F)) \\ & \text{gdw.} \quad \delta'(\pi(q_0), w) \in F' \\ & \text{gdw.} \quad \delta'(q'_0, w) \in F' \quad (\text{wegen } q'_0 = \pi(q_0)) \\ & \text{gdw.} \quad w \in L(\mathcal{A}') \end{aligned} \quad \square$$

Wir können nun Minimalität und Eindeutigkeit des reduzierten Automaten zeigen.

Satz 2.22

Es sei L eine erkennbare Sprache. Dann gilt:

- 1) \mathcal{A}_L hat minimale Zustandszahl unter allen DEAs, welche L erkennen.
- 2) Ist \mathcal{A} ein DEA mit $L(\mathcal{A}) = L$, so ist der reduzierte Automat $\mathcal{A}_{red} := \tilde{\mathcal{A}}_0$ isomorph zu \mathcal{A}_L .

Beweis.

- 1) Mit Satz 2.18 hat \simeq_L endlichen Index. Daher ist mit Lemma 2.17 \mathcal{A}_L ein DEA für L . Mit Lemma 2.15, 4) hat jeder DEA, der L erkennt, mindestens soviele Zustände, wie \simeq_L Klassen (d.h. wie \mathcal{A}_L Zustände) hat.
- 2) Es sei $\mathcal{A}_{red} = (Q, \Sigma, q_0, \delta, F)$ und $\mathcal{A}_L = (Q', \Sigma, q'_0, \delta', F')$. Wir definieren eine Funktion $\pi : Q \rightarrow Q'$ und zeigen, dass sie ein Isomorphismus ist. Für $q \in Q$ existiert (mindestens) ein Wort $w_q \in \Sigma^*$ mit $\delta(q_0, w_q) = q$, da in \mathcal{A}_{red} alle Zustände erreichbar sind. O.B.d.A. sei $w_{q_0} = \varepsilon$. Wir definieren $\pi(q) := [w_q]_L$.

I) π ist injektiv:

Wir müssen zeigen, dass aus $p \neq q$ auch $[w_p]_L \neq [w_q]_L$ folgt.

Da \mathcal{A}_{red} reduziert ist, sind verschiedene Zustände nicht äquivalent. Es gibt also mindestens ein w , für das

$$\delta(p, w) \in F \Leftrightarrow \delta(q, w) \in F$$

nicht gilt.

Das heißt aber, dass

$$\delta(q_0, w_p w) \in F \Leftrightarrow \delta(q_0, w_q w) \in F$$

nicht gilt und damit wiederum, dass $w_p w \in L \Leftrightarrow w_q w \in L$ nicht gilt. Also ist $w_p \not\equiv_L w_q$, d.h. $[w_p]_L \neq [w_q]_L$.

II) π ist surjektiv:

Folgt aus Injektivität und $|Q| \geq |Q'|$ (Aussage 1) des Lemmas).

III) $\pi(q_0) = q'_0$:

Da $w_{q_0} = \varepsilon$ und $q'_0 = [\varepsilon]_L$.

IV) $\pi(F) = F'$:

$q \in F \quad \text{gdw.} \quad \delta(q_0, w_q) = q \in F \quad \text{gdw.} \quad w_q \in L \quad \text{gdw.} \quad [w_q]_L \in F'$

V) $\pi(\delta(q, a)) = \delta'(\pi(q), a)$:

Es sei $\delta(q, a) =: p$. Dann ist $\pi(\delta(q, a)) = [w_p]_L$.

Außerdem ist $\delta'(\pi(q), a) = \delta'([w_q]_L, a) = [w_q a]_L$.

Es bleibt also $[w_p]_L = [w_q a]_L$ zu zeigen, d.h. $\forall w : w_p w \in L \quad \text{gdw.} \quad w_q a w \in L$.

Dies ist offenbar dann der Fall, wenn $\delta(q_0, w_p) = \delta(q_0, w_q a)$ ist:

$$\delta(q_0, w_p a) = \delta(\delta(q_0, w_q), a) = \delta(q, a) = p = \delta(q_0, w_p)$$

□

Im Prinzip liefert dieser Satz eine Methode, um von zwei Automaten zu entscheiden, ob sie dieselbe Sprache akzeptieren:

Korollar 2.23

Es seien \mathcal{A} und \mathcal{A}' DEAs. Dann gilt: $L(\mathcal{A}) = L(\mathcal{A}') \quad \text{gdw.} \quad \mathcal{A}_{red} \simeq \mathcal{A}'_{red}$.

Man kann die reduzierten Automaten wie beschrieben konstruieren. Für gegebene Automaten kann man feststellen, ob sie isomorph sind (teste alle Bijektionen). Hat man NEAs gegeben, so kann man diese zuerst deterministisch machen und dann das Korollar anwenden.

3. Nachweis der Nichterkennbarkeit

Um nachzuweisen, dass eine gegebene Sprache erkennbar ist, genügt es, einen endlichen Automaten (DEA oder NEA) dafür anzugeben. Der Nachweis, dass eine Sprache nicht erkennbar ist, gestaltet sich schwieriger: Es genügt ja nicht zu sagen, dass man keinen Automaten dafür gefunden hat. (Es gibt unendlich viele Automaten).

Wir haben bereits gesehen, dass man den Satz von Nerode (Satz 2.18) dazu verwenden kann. Ein anderes Hilfsmittel hierfür ist das Pumping-Lemma:

Lemma 3.1 (Pumping-Lemma, einfache Version)

Es sei L eine erkennbare Sprache. Dann gibt es eine natürliche Zahl $n_0 \geq 1$, so dass gilt: Jedes Wort $w \in L$ mit $|w| \geq n_0$ lässt sich zerlegen in $w = xyz$ mit

- $y \neq \varepsilon$
- $xy^kz \in L$ für alle $k \geq 0$.

Beweis. Es sei $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ ein NEA mit $L(\mathcal{A}) = L$. Wir wählen $n_0 = |Q|$. Für jedes Wort $w = a_1 \dots a_m \in L$ existiert ein Pfad $(p_0, a_1, p_1)(p_1, a_2, p_2) \dots (p_{m-1}, a_m, p_m)$ in \mathcal{A} mit $p_0 = q_0$ und $p_m \in F$.

Ist $m \geq n_0$, so können die $m + 1$ Zustände p_0, \dots, p_m nicht alle verschieden sein, da $|Q| = n_0 < m + 1$. Es gibt daher ein $i < j$ mit $p_i = p_j$.

Für $x := a_1 \dots a_i$, $y := a_{i+1} \dots a_j$, $z := a_{j+1} \dots a_m$ gilt daher $y \neq \varepsilon$ (da $i < j$) und $q_0 = p_0 \xrightarrow{x} p_i \xrightarrow{y} p_j \xrightarrow{z} p_m \in F$.

Folglich gilt für alle $k \geq 0$ auch $p_i \xrightarrow{y^k} p_i$, was $xy^kz \in L$ zeigt. □

Wir zeigen mit Hilfe dieses Lemmas (nochmals), dass die Sprache $\{a^n b^n \mid n \geq 0\}$ (vgl. Beispiel 2.19) nicht erkennbar ist.

Beispiel:

$L = \{a^n b^n \mid n \geq 0\}$ ist *nicht* erkennbar.

Beweis. Angenommen, L ist erkennbar. Dann gibt es eine Zahl n_0 mit den in Lemma 3.1 beschriebenen Eigenschaften. Es sei nun n so, dass $2n \geq n_0$ ist. Da $a^n b^n \in L$ ist, hat es eine Zerlegung $a^n b^n = xyz$ mit $|y| \geq 1$ und $xy^kz \in L$ für alle $k \geq 0$.

1. Fall: y liegt ganz in a^n .

D.h. $x = a^{n_1}$, $y = a^{n_2}$, $z = a^{n_3} b^n$ mit $n_2 > 0$ und $n = n_1 + n_2 + n_3$. Damit ist aber $xz = xy^0z = a^{n_1+n_3} b^n \notin L$, da $n_1 + n_3 < n$. Widerspruch.

2. Fall: y liegt ganz in b^n .

Führt entsprechend zu einem Widerspruch.

3. Fall: y enthält as und bs .

D.h. $x = a^{n_1}$, $y = a^{n_2} b^{n'_1}$, $z = b^{n'_2}$ mit $n_2 \neq 0 \neq n'_1$. Dann ist aber

$$xyyz = a^{n_1} a^{n_2} b^{n'_1} a^{n_2} b^{n'_1} b^{n'_2} \notin L$$

(da nach bs nochmal as kommen). Widerspruch.

In allen drei Fällen haben wir also einen Widerspruch erhalten, d.h. die Annahme „ L ist erkennbar“ war falsch. \square

Als eine weitere Konsequenz von Lemma 3.1 erhält man, dass das Leerheitsproblem für erkennbare Sprachen entscheidbar ist.

Satz 3.2

Es sei L eine erkennbare Sprache (gegeben durch NEA oder DEA). Dann kann man effektiv entscheiden, ob $L = \emptyset$ oder nicht.

Beweis. Es sei L erkennbar und n_0 die zugehörige Zahl aus Lemma 3.1. Dann gilt:

$$(\star) \quad L \neq \emptyset \text{ gdw. } \exists w \in L \text{ mit } |w| < n_0$$

Um $L = \emptyset$ zu entscheiden, muss man also nur für die endlich vielen Wörter $w \in \Sigma^*$ der Länge $< n_0$ entscheiden, ob w zu L gehört. Dies kann man für jedes einzelne Wort (z.B. durch Eingabe in den zugehörigen DEA) einfach entscheiden (vgl. Wortproblem, Satz 4.4).

Beweis von (\star) :

„ \Leftarrow “: trivial

„ \Rightarrow “: Es sei $L \neq \emptyset$. Wähle ein Wort w kürzester Länge in L . Wäre $|w| \geq n_0$, so müsste es eine Zerlegung $w = xyz$, $y \neq \varepsilon$ geben mit $xy^0z = xz \in L$. Dies ist ein Widerspruch zur Minimalität von w . \square

Mit Hilfe der einfachen Variante des Pumping-Lemmas gelingt es nicht immer, Nichterkennbarkeit nachzuweisen.

Beispiel 3.3

Ist $L = \{a^n b^m \mid n \neq m\}$ erkennbar? Versucht man, Nichterkennbarkeit mit Lemma 3.1 zu zeigen, so scheitert man, da das Lemma für L zutrifft:

Wähle $n_0 := 3$. Es sei nun $w \in L$ mit $|w| \geq 3$, d.h. $w = a^n b^m$, $n \neq m$ und $n + m \geq 3$.

1. Fall: $n > m$

D.h. $n = m + i$ für $i \geq 1$.

1.1.: $i > 1$, d.h. $n - 1 > m$.

Für $x = \varepsilon$, $y = a$, $z = a^{n-1}b^m$ gilt für alle $k \geq 0$: $xy^kz = a^k a^{n-1} b^m \in L$, da $k + n - 1 \geq n - 1 > m$ (wegen $n - m = i > 1$).

1.2.: $i = 1$, d.h. $n = m + 1$.

Wegen $n + m \geq 3$ folgt $n = m + 1 \geq 2$.

Für $x = \varepsilon$, $y = a^2$, $z = a^{n-2}b^m$ gilt

a) $xy^0z \in L$, da $n - 2 = m - 1 \neq m$

b) $xy^kz \in L$ für $k \geq 1$, da $(n - 2) + k \cdot 2 \geq n > m$

2. Fall: $n < m$

kann entsprechend behandelt werden.

Trotzdem ist $L = \{a^n b^m \mid n \neq m\}$ nicht erkennbar, was man mit der folgenden verschärften Variante des Pumping-Lemmas nachweisen kann.

Lemma 3.4 (Pumping-Lemma, verschärfte Variante)

Es sei L erkennbar. Dann gibt es eine natürliche Zahl $n_0 \geq 1$, so dass gilt:

Für alle Wörter u, v, w mit $uvw \in L$ und $|v| \geq n_0$ gibt es eine Zerlegung $v = xyz$ mit

- $y \neq \varepsilon$
- $uxy^kzw \in L$ für alle $k \geq 0$

Beweis. Es sei wieder $n_0 := |Q|$, wobei Q die Zustände eines NEA \mathcal{A} für L sind. Ist $uvw \in L$, so gibt es Zustände $p, q, f \in Q$ mit

$$q_0 \xrightarrow{u}_{\mathcal{A}} p \xrightarrow{v}_{\mathcal{A}} q \xrightarrow{w}_{\mathcal{A}} f \in F$$

Auf dem Pfad von p nach q liegen $|v| + 1 > n_0$ Zustände, also müssen zwei davon gleich sein. Jetzt kann man wie im Beweis von Lemma 3.1 weitermachen. □

Was ist der Vorteil dieses Lemmas beim Nachweis der Nichterkennbarkeit? Man weiß, dass man bereits beliebig lange Teilstücke zerlegen kann. Dadurch kann man das y geeignet positionieren (im Beispiel 3.3 im kürzeren Block).

Beispiel 3.3 (Fortsetzung)

$L = \{a^n b^m \mid n \neq m\}$ ist *nicht* erkennbar.

Beweis. Angenommen, L ist doch erkennbar; dann gibt es $n_0 \geq 1$, das die in Lemma 3.4 geforderten Eigenschaften hat. Wir betrachten nun die Wörter

$$u := \varepsilon, \quad v := a^{n_0}, \quad w := b^{n_0! + n_0}$$

Offenbar ist $uvw = a^{n_0} b^{n_0! + n_0} \in L$. Mit Lemma 3.4 gibt es eine Zerlegung $v = xyz$ mit $y \neq \varepsilon$ und $uxy^kzw \in L$ für alle $k \geq 0$. Es sei

$$x = a^{n_1}, \quad y = a^{n_2}, \quad z = a^{n_3}, \quad n_1 + n_2 + n_3 = n_0, \quad n_2 > 0$$

Offenbar existiert ein l mit $n_2 \cdot l = n_0!$ (da $0 < n_2 \leq n_0$). Es ist nun aber

$$n_1 + (l + 1) \cdot n_2 + n_3 = n_0 + n_0!$$

was $uxy^{l+1}zw \notin L$ liefert (Widerspruch). □

4. Abschlusseigenschaften und Entscheidungsprobleme

Wir zeigen zunächst, dass die Klasse der erkennbaren Sprachen unter den Booleschen Operationen sowie Konkatenation und Kleene-Stern abgeschlossen ist.

Satz 4.1 (Abschluss erkennbarer Sprachen)

Sind L_1 und L_2 erkennbar, so sind auch

- $L_1 \cup L_2$ (Vereinigung)
- $\overline{L_1}$ (Komplement)
- $L_1 \cap L_2$ (Durchschnitt)
- $L_1 \setminus L_2$ (Differenz)
- $L_1 \cdot L_2$ (Konkatenation)
- L_1^* (Kleene-Stern)

erkennbar.

Beweis. Es seien $\mathcal{A}_i = (Q_i, \Sigma, q_{0i}, \Delta_i, F_i)$ zwei NEAs für L_i ($i = 1, 2$). O.B.d.A. gelte $Q_1 \cap Q_2 = \emptyset$.

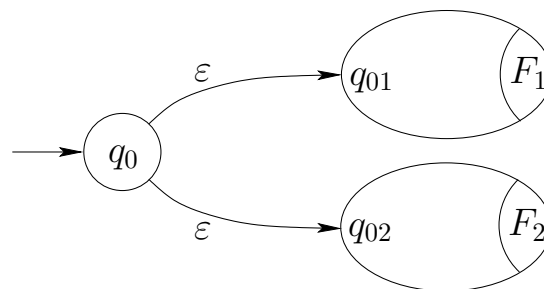
1) Abschluss unter Vereinigung:

Der folgende ε -NEA erkennt $L_1 \cup L_2$:

$\mathcal{A} := (Q_1 \cup Q_2 \cup \{q_0\}, \Sigma, q_0, \Delta, F_1 \cup F_2)$, wobei

- $q_0 \notin Q_1 \cup Q_2$ und
- $\Delta := \Delta_1 \cup \Delta_2 \cup \{(q_0, \varepsilon, q_{01}), (q_0, \varepsilon, q_{02})\}$.

Schematisch sieht der Vereinigungsautomat \mathcal{A} so aus.



Mit Lemma 1.12 gibt es zu \mathcal{A} einen äquivalenten NEA.

2) Abschluss unter Komplement:

Einen DEA für $\overline{L_1}$ erhält man wie folgt:

Zunächst verwendet man die Potenzmengenkonstruktion, um zu \mathcal{A}_1 einen äquivalenten DEA $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ zu konstruieren. Der DEA für $\overline{L_1}$ ist nun $\overline{\mathcal{A}} := (Q, \Sigma, q_0, \delta, Q \setminus F)$. Es gilt nämlich:

$$\begin{aligned}
 w \in \overline{L_1} & \quad \text{gdw.} \quad w \notin L(\mathcal{A}_1) \\
 & \quad \text{gdw.} \quad w \notin L(\mathcal{A}) \\
 & \quad \text{gdw.} \quad \delta(q_0, w) \notin F \\
 & \quad \text{gdw.} \quad \delta(q_0, w) \in Q \setminus F \\
 & \quad \text{gdw.} \quad w \in L(\overline{\mathcal{A}})
 \end{aligned}$$

Beachte: Man darf dies nicht direkt mit dem NEA machen!

3) **Abschluss unter Durchschnitt:**

Durch Ausnutzen von $L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$ folgt 3) aus 1) und 2).

Da die Potenzmengenkonstruktion aber sehr aufwendig sein kann, ist es günstiger, direkt einen NEA für $L_1 \cap L_2$ zu konstruieren, den sogenannten *Produktautomaten*:

$$\mathcal{A} := (Q_1 \times Q_2, \Sigma, (q_{01}, q_{02}), \Delta, F_1 \times F_2)$$

mit

$$\Delta := \{((q_1, q_2), a, (q'_1, q'_2)) \mid (q_1, a, q'_1) \in \Delta_1 \text{ und } (q_2, a, q'_2) \in \Delta_2\}$$

Ein Übergang in \mathcal{A} ist also genau dann möglich, wenn der entsprechende Übergang in \mathcal{A}_1 und \mathcal{A}_2 möglich ist. Daraus ergibt sich leicht $L(\mathcal{A}) = L_1 \cap L_2$.

4) **Abschluss unter Differenz:**

Folgt aus 1) und 2), da

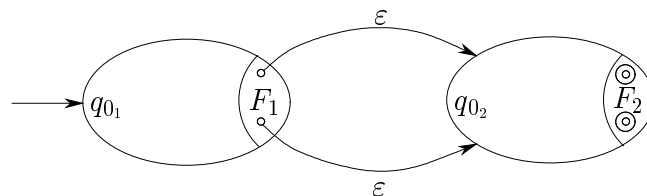
$$L_1 \setminus L_2 = L_1 \cap \overline{L_2}.$$

5) **Abschluss unter Konkatenation:**

Der folgende ε -NEA erkennt $L_1 \cdot L_2$:

$\mathcal{A} := (Q_1 \cup Q_2, \Sigma, q_{01}, \Delta, F_2)$, wobei

$$\Delta := \Delta_1 \cup \Delta_2 \cup \{(f, \varepsilon, q_{02}) \mid f \in F_1\}$$

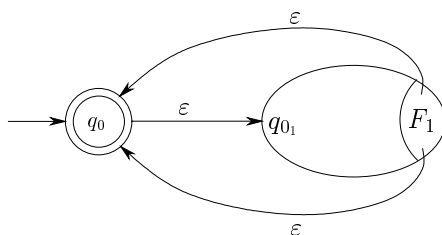


6) **Abschluss unter Kleene-Stern:**

Der folgende ε -NEA erkennt L_1^* :

$\mathcal{A} := (Q_1 \cup \{q_0\}, \Sigma, q_0, \Delta, \{q_0\})$, wobei

- $q_0 \notin Q_1$
- $\Delta := \Delta_1 \cup \{(f, \varepsilon, q_0) \mid f \in F_1\} \cup \{(q_0, \varepsilon, q_{01})\}$.



□

Beachte:

Alle angegebenen Konstruktionen sind effektiv. Die Automaten für die Sprachen $L_1 \cup L_2$, $L_1 \cap L_2$, $L_1 \cdot L_2$ und L_1^* sind polynomiell in der Größe der Automaten für L_1 , L_2 . Beim Komplement kann die Konstruktion exponentiell sein, wenn man von einem NEA ausgeht.

Man kann derartige Abschlusseigenschaften auch dazu verwenden, Nichterkennbarkeit einer Sprache L nachzuweisen.

Beispiel 4.2

$L := \{a^n b^m \mid n \neq m\}$ ist *nicht* erkennbar (vgl. Beispiel 3.3). Anstatt dies direkt mit Lemma 3.4 zu zeigen, kann man auch verwenden, dass bereits bekannt ist, dass die Sprache $L' := \{a^n b^n \mid n \geq 0\}$ *nicht* erkennbar ist. Wäre nämlich L erkennbar, so auch $L' = \bar{L} \cap \{a\}^* \cdot \{b\}^*$. Da wir schon wissen, dass L' nicht erkennbar ist, kann auch L nicht erkennbar sein.

Wir betrachten nun drei Probleme im Zusammenhang mit erkennbaren Sprachen, für die wir Entscheidbarkeit und Komplexität untersuchen, und zwar:

- das *Leerheitsproblem*
- das *Wortproblem* und
- das *Äquivalenzproblem*.

Dabei gehen wir davon aus, dass die erkennbare Sprache durch einen DEA oder NEA gegeben ist. Bezüglich der Entscheidbarkeit dieser Probleme macht es keinen Unterschied, ob man einen DEA oder einen NEA gegeben hat, da man aus einem NEA effektiv einen äquivalenten DEA konstruieren kann (Satz 2.4). Bezüglich der Komplexität kann der Unterschied aber relevant sein, da der Übergang NEA \rightarrow DEA exponentiell sein kann.

Leerheitsproblem:

Geg.: erkennbare Sprache L (durch DEA oder NEA)

Frage: Ist $L \neq \emptyset$?

Wir wissen bereits (Satz 3.2), dass das Problem entscheidbar ist. Allerdings ist das im Beweis des Satzes beschriebene Entscheidungsverfahren viel zu aufwendig (*exponentiell* viele Wörter der Länge $< n_0$, falls $|\Sigma| > 1$).

Satz 4.3

Es sei $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ ein NEA. Dann kann man das Problem „ $L(\mathcal{A}) \neq \emptyset$?“ in der Zeit $O(|Q| + |\Delta|)$ entscheiden.

Beweis. Man kann \mathcal{A} als gerichteten Graphen $G = (Q, E)$ auffassen mit

$$E := \{(q_1, q_2) \mid (q_1, a, q_2) \in \Delta \text{ für ein } a \in \Sigma\}$$

Dann gilt: $L(\mathcal{A}) \neq \emptyset$ gdw. in der von q_0 aus erreichbaren Knotenmenge befindet sich ein Endzustand. Die von q_0 aus erreichbaren Knoten kann man mit Aufwand $O(|Q| + |E|)$ berechnen (vgl. Vorlesung „Algorithmen und Datenstrukturen“). Insbesondere ist also das Leerheitsproblem für erkennbare Sprachen mit polynomiellem Aufwand lösbar. \square

Wortproblem:

Geg.: erkennbare Sprache L , Wort $w \in \Sigma^*$

Frage: Gilt $w \in L$?

Ist $L = L(\mathcal{A})$ für einen DEA $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$, so kann man einfach, beginnend mit q_0 , durch Anwendung von δ berechnen, zu welchem Zustand in \mathcal{A} man mit w kommt und prüfen, ob dies ein Endzustand ist. Nimmt man Σ als konstant an, so benötigt eine Anwendung von δ nur konstante Zeit. Dies liefert:

Satz 4.4

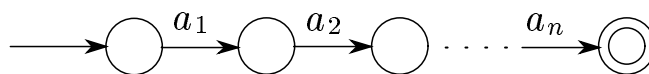
Es sei $\mathcal{A} = (Q, \Sigma, q_0, \delta, F)$ ein DEA und $w \in \Sigma^$. Dann kann man „ $w \in L(\mathcal{A})$?“ in der Zeit $O(|w|)$ entscheiden.*

Für einen NEA ist dies nicht so einfach, da man ja verschiedene mit w beschriftete Pfade haben kann und man diese (im schlimmsten Fall) alle betrachten muss, um festzustellen, ob einer davon mit einem Endzustand aufhört.

Satz 4.5

Es sei $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ ein NEA und $w \in \Sigma^$. Dann kann man „ $w \in L(\mathcal{A})$?“ in der Zeit $O(|w| \cdot (|Q| + |\Delta|))$ entscheiden.*

Beweis. Konstruiere zunächst einen Automaten \mathcal{A}_w , der genau das Wort $w = a_1 \dots a_n$ akzeptiert:



Dieser Automat hat $|w| + 1$ Zustände. Offenbar ist $w \in L(\mathcal{A})$ gdw. $L(\mathcal{A}) \cap L(\mathcal{A}_w) \neq \emptyset$. Wie groß ist nun der Produktautomat zu \mathcal{A} und \mathcal{A}_w ?

Zustände: $|Q| \cdot (|w| + 1)$

Übergänge: Für jeden Übergang $\bigcirc \xrightarrow{a_i} \bigcirc$ gibt es maximal $|\Delta|$ mögliche Übergänge in \mathcal{A} .

Also maximal $|w| \cdot |\Delta|$ viele Übergänge im Produktautomaten.

Nach Satz 4.3 ist daher der Aufwand zum Testen von $L(\mathcal{A}) \cap L(\mathcal{A}_w) \neq \emptyset$:

$$O(|Q| \cdot (|w| + 1) + |w| \cdot |\Delta|) = O(|w| \cdot (|Q| + |\Delta|)) \quad \square$$

Äquivalenzproblem:

Geg.: erkennbare Sprachen L_1, L_2

Frage: Gilt $L_1 = L_2$?

Wie bereits erwähnt, kann man das Äquivalenzproblem auf das Leerheitsproblem *reduzieren*:

$$L_1 = L_2 \quad \text{gdw.} \quad (L_1 \cap \overline{L_2}) \cup (\overline{L_1} \cap L_2) = \emptyset$$

Satz 4.6

Das Äquivalenzproblem ist für erkennbare Sprachen entscheidbar. Sind die Sprachen durch DEAs gegeben, so ist dies in polynomieller Zeit möglich.

Beweis. Wir haben gesehen, dass die Automatenkonstruktionen, welche Abschluss unter Vereinigung, Durchschnitt und Komplement zeigen, effektiv sind. Daraus ergibt sich direkt die Entscheidbarkeit des Äquivalenzproblems (auch bei NEAs). Die Konstruktionen für Vereinigung und Durchschnitt sind polynomiell. Beim Komplement ist dies nur dann der Fall, wenn bereits DEAs vorliegen. \square

Wir werden später sehen, dass sich der exponentielle Zeitaufwand für die Potenzmengenkonstruktion bei NEAs (wahrscheinlich) nicht vermeiden lässt: das Äquivalenzproblem für NEAs ist *PSPACE*-vollständig (schlimmer als *NP*).

5. Reguläre Ausdrücke und Sprachen

Wir haben bisher *verschiedene äquivalente Charakterisierungen* der Klasse der erkennbaren Sprachen mit Hilfe von Transitionssystemen und Rechtskongruenzen gesehen:

Eine Sprache $L \subseteq \Sigma^*$ ist *erkennbar* gdw.

- (1) $L = L(\mathcal{A})$ für einen NEA \mathcal{A} .
- (2) $L = L(\mathcal{A})$ für einen ε -NEA \mathcal{A} .
- (3) $L = L(\mathcal{A})$ für einen NEA mit Wortübergängen \mathcal{A} .
- (4) $L = L(\mathcal{A})$ für ein endliches Transitionssystem \mathcal{A} .
- (5) $L = L(\mathcal{A})$ für einen DEA \mathcal{A} .
- (6) Die Nerode-Rechtskongruenz \simeq_L hat endlichen Index.

Im folgenden betrachten wir eine weitere Charakterisierung mit Hilfe *regulärer Ausdrücke*.

Definition 5.1 (Syntax regulärer Ausdrücke)

Es sei Σ ein endliches Alphabet. Die Menge Reg_Σ der *regulären Ausdrücke über Σ* ist induktiv definiert:

- $\emptyset, \varepsilon, a$ (für $a \in \Sigma$) sind Elemente von Reg_Σ .
- Sind $r, s \in Reg_\Sigma$, so auch $(r + s), (r \cdot s), r^* \in Reg_\Sigma$.

Beispiel 5.2

$((a \cdot b^*) + \emptyset^*)^* \in Reg_\Sigma$ für $\Sigma = \{a, b\}$

Notation:

Um Klammern zu sparen, lassen wir Außenklammern weg und vereinbaren,

- dass $*$ stärker bindet als \cdot
- dass \cdot stärker bindet als $+$
- \cdot lassen wir meist ganz wegfallen.

Der Ausdruck aus Beispiel 5.2 kann also geschrieben werden als $(ab^* + \emptyset^*)^*$.

Jedem regulären Ausdruck r über Σ wird eine formale Sprache $L(r)$ zugeordnet.

Definition 5.3 (Semantik regulärer Ausdrücke)

Die durch den regulären Ausdruck r definierte Sprache $L(r)$ ist induktiv definiert:

- $L(\emptyset) := \emptyset, \quad L(\varepsilon) := \{\varepsilon\}, \quad L(a) := \{a\}$
- $L(r + s) := L(r) \cup L(s), \quad L(r \cdot s) := L(r) \cdot L(s), \quad L(r^*) := L(r)^*$

Eine Sprache $L \subseteq \Sigma^*$ heißt *regulär*, falls es ein $r \in Reg_\Sigma$ gibt mit $L = L(r)$.

Beispiel:

- $(a+b)^*ab(a+b)^*$ definiert die Sprache aller Wörter über $\{a, b\}$, die Infix ab haben.
- $L(ab^* + b) = \{ab^i \mid i \geq 0\} \cup \{b\}$

Bemerkung:

Statt $L(r)$ schreiben wir im folgenden häufig einfach r .

Dies ermöglicht es zu schreiben:

- $abb \in ab^* + b$ (eigentlich $abb \in L(ab^* + b)$)
- $(ab)^*a = a(ba)^*$ (eigentlich $L((ab)^*a) = L(a(ba)^*)$)

Wir zeigen nun, dass man mit regulären Ausdrücken genau die erkennbaren Sprachen definieren kann.

Satz 5.4 (Kleene)

Für eine Sprache $L \subseteq \Sigma^*$ sind äquivalent:

- 1) L ist regulär.
- 2) L ist erkennbar.

Beweis.

„1 \rightarrow 2“: Induktion über den Aufbau regulärer Ausdrücke

Verankerung:

- $L(\emptyset) = \emptyset$ erkennbar: $\rightarrow \bigcirc$ ist NEA für \emptyset (kein Endzustand).
- $L(\varepsilon) = \{\varepsilon\}$ erkennbar: $\rightarrow \bigcirc$ ist NEA für $\{\varepsilon\}$.
- $L(a) = \{a\}$ erkennbar: $\rightarrow \bigcirc \xrightarrow{a} \bigcirc$ ist NEA für $\{a\}$.

Schritt: Weiß man bereits, dass $L(r)$ und $L(s)$ erkennbar sind, so folgt mit Satz 4.1 (Abschlusseigenschaften), dass auch

- $L(r + s) = L(r) \cup L(s)$
- $L(r \cdot s) = L(r) \cdot L(s)$ und
- $L(r^*) = L(r)^*$

erkennbar sind.

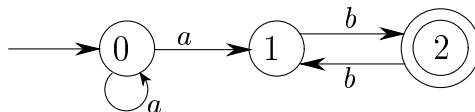
„2 \rightarrow 1“: Sei $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$ ein NEA mit $L = L(\mathcal{A})$. Wie in Definition 2.8 sei für $q \in Q$ der NEA $\mathcal{A}_q := (Q, \Sigma, q, \Delta, F)$ und $L_q = L(\mathcal{A}_q)$ definiert, d.h. insbesondere $L = L_{q_0}$.

Der Zusammenhang zwischen den L_q s kann nun als *Gleichungssystem* beschrieben werden, dessen Lösungen eindeutig bestimmte *reguläre Sprachen* sind.

Diesen Zusammenhang werden wir im Folgenden näher beleuchten.

Beispiel 5.5

Gegeben sei der folgende NEA \mathcal{A} .



Die Sprachen L_q für $q \in Q = \{0, 1, 2\}$ kann man bestimmen als

$$L_0 = \{a\} \cdot L_0 \cup \{a\} \cdot L_1$$

$$L_1 = \{b\} \cdot L_2$$

$$L_2 = \{b\} \cdot L_1 \cup \{\varepsilon\}$$

Allgemein:

Seien $p, q \in Q$. Wir definieren:

- $A_{p,q} = \{a \in \Sigma \mid (p, a, q) \in \Delta\}$ und
- $B_p = \begin{cases} \{\varepsilon\} & \text{falls } p \in F \\ \emptyset & \text{falls } p \notin F \end{cases}$

Damit erfüllen die Sprachen L_p die folgenden Gleichungen. Für alle $p \in Q$ gilt:

$$L_p = \left(\bigcup_{q \in Q} A_{p,q} \cdot L_q \right) \cup B_p$$

Behauptung:

Das Gleichungssystem

$$(\star) \quad X_p = \left(\bigcup_{q \in Q} A_{p,q} \cdot X_q \right) \cup B_p \quad (p \in Q)$$

hat *genau eine* Lösung und diese besteht aus *regulären* Sprachen.

Aus der Behauptung folgt, dass die Sprachen $X_p = L_p$ regulär sind, da sie ja das System (\star) lösen. Damit erhalten wir den Abschluss des Beweises von Satz 5.4, Richtung „ $2 \rightarrow 1$ “.

Wir zeigen zunächst, wie man eine einzige Gleichung der Form

$$(\star\star) \quad X = A \cdot X \cup B$$

lösen kann. Daraus ergibt sich dann die Lösung von Gleichungssystemen der Form (\star) durch Induktion über die Anzahl der Variablen.

Lemma 5.6 (Arden)

Es seien $A, B \subseteq \Sigma^*$ und $\varepsilon \notin A$. Die Gleichung $(\star\star) \quad X = A \cdot X \cup B$ hat als eindeutige Lösung $X = A^* \cdot B$.

Beachte:

Sind A, B regulär, so auch A^*B .

Aus dem Lemma folgt dann für die Gleichung (\star) , die wir als Gleichung der Form $(\star\star)$ für ein fixiertes $p \in Q$ auffassen:

$$X_p = A_{p,p} \cdot X_p \cup \left(\left(\bigcup_{p \neq q} A_{p,q} \cdot X_q \right) \cup B_p \right)$$

hat als eindeutige Lösung

$$X_p = A_{p,p}^* \cdot \left(\left(\bigcup_{q \neq p} A_{p,q} \cdot X_q \right) \cup B_p \right).$$

Setzt man diese Lösung in (\star) ein, so erhält man ein System mit einer Variablen weniger. Nach Induktion hat dieses eine eindeutige Lösung, die aus regulären Sprachen besteht. Da in obiger Lösung für X_p nur reguläre Operationen $(\cup, \cdot, ^*)$ verwendet werden, ist auch diese Lösung regulär (und eindeutig nach Arden).

Beispiel 5.5 (Fortsetzung)

$$\begin{aligned} X_0 &= \{a\} \cdot X_0 \cup \{a\} \cdot X_1 \\ X_1 &= \{b\} \cdot X_2 \\ X_2 &= \{b\} \cdot X_1 \cup \{\varepsilon\} \end{aligned}$$

Auflösen nach X_0 :

$$X_0 = \{a\}^* \cdot \{a\} \cdot X_1$$

Einsetzen ändert die anderen Gleichungen nicht.

Auflösen nach X_1 :

$$X_1 = \emptyset^* \cdot \{b\} \cdot X_2 = \{b\} \cdot X_2$$

Einsetzen liefert:

$$X_2 = \{b\} \cdot \{b\} \cdot X_2 \cup \{\varepsilon\}$$

Auflösen nach X_2 :

$$X_2 = (\{b\} \cdot \{b\})^*$$

Damit ist

$$X_1 = \{b\} \cdot (\{b\} \cdot \{b\})^*$$

und

$$X_0 = \{a\}^* \cdot \{a\} \cdot \{b\} \cdot (\{b\} \cdot \{b\})^* = L(\mathcal{A}).$$

Als regulären Ausdruck für $L(\mathcal{A})$ liefert dieses Verfahren also: $L(\mathcal{A}) = L(a^*ab(bb)^*)$.

Beweis von Lemma 5.6.

1) A^*B ist Lösung:

$$A \cdot A^* \cdot B \cup B = (A \cdot A^* \cup \{\varepsilon\}) \cdot B = A^* \cdot B$$

2) Eindeutigkeit:

Es sei L eine Lösung, d.h. $L = A \cdot L \cup B$.

2.1) $A^*B \subseteq L$.

Wir zeigen durch Induktion über n : $A^n B \subseteq L$.

- $A^0 B = B \subseteq (A \cdot L \cup B) = L$
- Gelte $A^n B \subseteq L$.
Es folgt: $A^{n+1} B = A \cdot A^n B \subseteq A \cdot L \subseteq A \cdot L \cup B = L$

Wegen $A^*B = \bigcup_{n \geq 0} A^n B$ folgt damit $A^*B \subseteq L$.

2.2) $L \subseteq A^*B$.

Angenommen, dies gilt nicht. Es sei $w \in L$ ein Wort minimaler Länge mit $w \notin A^*B$.

$w \in L = A \cdot L \cup B$, d.h. $w \in A \cdot L$ oder $w \in B$. Für $w \in B$ folgt $w \in A^*B$ (Widerspruch).

Aus $w \in A \cdot L$ folgt, dass es $u \in A$ und $v \in L$ gibt mit $w = uv$. Wegen $\varepsilon \notin A$ ist $|v| < |w|$. Wegen Minimalität von w folgt $v \in A^*B$ (Widerspruch, da dann $w = uv \in A^*B$).

□

Zum Abschluss von Teil I erwähnen wir einige hier aus Zeitgründen nicht behandelte Themenbereiche:

Endliche Automaten mit Ausgabe:

Übergänge $p \xrightarrow{a/v} q$, wobei $v \in \Gamma^*$ ein Wort über einem Ausgabealphabet ist. Solche Automaten beschreiben spezielle Funktionen $\Sigma^* \rightarrow \Gamma^*$.

algebraische Theorie formaler Sprachen:

Jeder Sprache L wird ein Monoid M_L (syntaktisches Monoid) zugeordnet. Klassen von Sprachen entsprechen dann Klassen von Monoiden, z.B. L ist regulär gdw. M_L ist endlich.

Automaten auf unendlichen Wörtern:

Büchi-Automaten sind endliche Automaten, bei denen man unendliche Wörter (unendliche Folgen von Symbolen) eingibt.

Baumautomaten:

Sind Automaten, die Bäume statt Wörter als Eingaben haben.

II. Grammatiken, kontextfreie Sprachen und Kellerautomaten

Einführung

Wir werden hier Klassen formaler Sprachen untersuchen, die allgemeiner sind als die der regulären Sprachen. Insbesondere werden wir die Klasse der kontextfreien Sprachen genauer betrachten, da durch sie z.B. die Syntax von Programmiersprachen (zumindest in großen Teilen) beschreibbar ist.

Zunächst führen wir allgemein den Begriff der Grammatik ein. Klassen formaler Sprachen erhält man durch einschränkende Bedingungen an die Form der Grammatik.

6. Die Chomsky-Hierarchie

Grammatiken dienen dazu, Wörter zu erzeugen. Man hat dazu *Regeln*, die es erlauben, ein Wort durch ein anderes Wort zu ersetzen (aus ihm abzuleiten). Die *erzeugte Sprache* ist die Menge der Wörter, die ausgehend von einem *Startsymbol* erzeugt werden können durch wiederholtes Ersetzen.

Beispiel 6.1

$$\begin{aligned} \text{Regeln:} \quad S &\longrightarrow aSb & (1) \\ S &\longrightarrow \varepsilon & (2) \end{aligned}$$

Startsymbol: S

Eine mögliche Ableitung eines Wortes wäre:

$$S \xrightarrow{1} aSb \xrightarrow{1} aaSbb \xrightarrow{1} aaaSbbb \xrightarrow{2} aaabbb$$

Das Symbol S ist hier ein Hilfssymbol (*nichtterminales Symbol*) und man ist nur an den erzeugten Wörtern interessiert, die das Hilfssymbol nicht enthalten (*Terminalwörter*). Man sieht leicht, dass dies hier alle Wörter $a^n b^n$ mit $n \geq 0$ sind.

Definition 6.2 (Grammatik)

Eine *Grammatik* ist von der Form $G = (N, \Sigma, P, S)$, wobei

- N und Σ endliche, disjunkte Alphabete sind. (Man bezeichnet die Symbole aus N als *Nichtterminalsymbole*, die Symbole aus Σ als *Terminalsymbole*),
- $S \in N$ das *Startsymbol* ist,
- $P \subseteq (N \cup \Sigma)^+ \times (N \cup \Sigma)^*$ eine endliche Menge von Ersetzungsregeln (*Produktionen*) ist.

Produktionen $(u, v) \in P$ schreibt man gewöhnlich als $u \longrightarrow v$.

Beispiel 6.3

$G = (N, \Sigma, P, S)$ mit

- $N = \{S, B\}$
- $\Sigma = \{a, b, c\}$
- $P = \{S \longrightarrow aSBc,$
 $S \longrightarrow abc,$
 $cB \longrightarrow Bc,$
 $bB \longrightarrow bb\}$

Im Folgenden schreiben wir meistens Elemente von N mit Grossbuchstaben und Elemente von Σ mit Kleinbuchstaben.

Wir definieren nun, was es heißt, dass man ein Wort durch Anwenden der Regeln aus einem anderen ableiten kann.

Definition 6.4 (durch eine Grammatik erzeugte Sprache)

Es sei $G = (N, \Sigma, P, S)$ eine Grammatik und x, y seien Wörter aus $(N \cup \Sigma)^*$.

- 1) y aus x direkt ableitbar:
 $x \vdash_G y$ gdw. $\exists x_1, x_2 \in (N \cup \Sigma)^* : \exists u \rightarrow v \in P : x = x_1 u x_2 \wedge y = x_1 v x_2$
- 2) y aus x in n Schritten ableitbar:
 $x \vdash_G^n y$ gdw. $\exists x_0, x_1, \dots, x_n \in (N \cup \Sigma)^* : x_0 = x \wedge x_n = y \wedge x_i \vdash_G x_{i+1}$ für $0 \leq i < n$
- 3) y aus x ableitbar:
 $x \vdash_G^* y$ gdw. $\exists n \geq 0 : x \vdash_G^n y$
- 4) Die durch G erzeugte Sprache ist
 $L(G) := \{w \in \Sigma^* \mid S \vdash_G^* w\}$.

Man ist also bei der erzeugten Sprache nur an den in G aus S ableitbaren Terminalwörtern interessiert.

Beispiel 6.3 (Fortsetzung)

$$\begin{aligned}
 S &\vdash_G abc, \text{ d.h. } abc \in L(G) \\
 S &\vdash_G aSBc \\
 &\vdash_G aaSBcBc \\
 &\vdash_G aaabcBcBc \\
 &\vdash_G aaabBccBc \\
 &\vdash_G^2 aaabBBccc \\
 &\vdash_G^2 aaabbbccc
 \end{aligned}$$

Es gilt: $L(G) = \{a^n b^n c^n \mid n \geq 1\}$

Beweis.

„ \supseteq “: Für $n = 1$ ist $abc \in L(G)$ klar. Für $n > 1$ sieht man leicht:

$$S \vdash_G^{n-1} a^{n-1} S (Bc)^{n-1} \vdash_G a^n bc (Bc)^{n-1} \vdash_G^* a^n b B^{n-1} c^n \vdash_G^* a^n b^n c^n$$

„ \subseteq “: Gelte $S \vdash_G^* w$ mit $w \in \Sigma^*$. Wird sofort die Produktion $S \rightarrow abc$ verwendet, so ist $w = abc \in \{a^n b^n c^n \mid n \geq 1\}$.

Sonst betrachten wir die Stelle in der Ableitung, an der S das letzte mal auftritt:

$$S \vdash_G^* a^{n-1} S u \text{ mit } u \in \{c, B\}^* \text{ und } |u|_B = |u|_c = n - 1$$

(nur $S \rightarrow aSBc, cB \rightarrow Bc$ angewendet, da noch kein b vorhanden).

$$a^{n-1} S u \vdash_G a^n bc u \vdash_G^* w.$$

Um von $a^n bc u$ aus nun alle nichtterminalen Symbole B zu entfernen, muss man $bB \rightarrow bb$ anwenden. Damit B aber auf b stößt, muss es links von allen c s stehen.

□

Beispiel 6.5

$G = (N, \Sigma, P, S)$ mit

- $N = \{S, B\}$
- $\Sigma = \{a, b\}$
- $P = \{S \rightarrow aS,$
 $S \rightarrow bS,$
 $S \rightarrow abB,$
 $B \rightarrow aB,$
 $B \rightarrow bB,$
 $B \rightarrow \varepsilon\}$

$$L(G) = \Sigma^* \cdot \{a\} \cdot \{b\} \cdot \Sigma^*$$

Die Grammatiken aus Beispiel 6.5, 6.3 und 6.1 gehören zu unterschiedlichen Stufen der **Chomsky-Hierarchie**:

Definition 6.6 (Typen von Chomsky-Grammatiken)

Es sei $G = (N, \Sigma, P, S)$ eine Grammatik.

- Jede Grammatik G heißt Grammatik vom **Typ 0**.
- G heißt Grammatik vom **Typ 1 (kontextsensitiv)**, falls jede Produktion von G die Form
 - $u_1Au_2 \rightarrow u_1wu_2$ mit $A \in N, u_1, u_2, w \in (\Sigma \cup N)^*$ und $|w| \geq 1$ oder
 - $S \rightarrow \varepsilon$ hat.

Ist $S \rightarrow \varepsilon \in P$, so kommt S nicht auf der rechten Seite einer Produktion vor.

- G heißt Grammatik vom **Typ 2 (kontextfrei)**, falls jede Regel von G die Form $A \rightarrow w$ hat mit $A \in N, w \in (\Sigma \cup N)^*$.
- G heißt Grammatik vom **Typ 3 (rechtslinear)**, falls jede Regel von G die Form $A \rightarrow uB$ oder $A \rightarrow u$ hat mit $A, B \in N, u \in \Sigma^*$.

kontextfrei: Die linke Seite jeder Produktion besteht nur aus einem Nichtterminalsymbol A , das daher stets unabhängig vom Kontext im Wort ersetzt werden kann.

kontextsensitiv: $u_1Au_2 \rightarrow u_1wu_2$. Hier ist die Ersetzung von A durch w abhängig davon, dass der richtige Kontext (u_1 links und u_2 rechts) vorhanden ist. $|w| \geq 1$ sorgt dafür, dass die Produktionen die Wörter verlängern (Ausnahme $S \rightarrow \varepsilon$, die aber nur zur Erzeugung von ε dient).

Definition 6.7 (Klasse der Typ- i -Sprachen)

Für $i = 0, 1, 2, 3$ ist die *Klasse der Typ- i -Sprachen* definiert als $\mathcal{L}_i := \{L(G) \mid G \text{ ist Grammatik vom Typ } i\}$.

Wir werden sehen: $\mathcal{L}_3 \subset \mathcal{L}_2 \subset \mathcal{L}_1 \subset \mathcal{L}_0$.

Nach Definition der Grammatiktypen gilt offenbar $\mathcal{L}_3 \subseteq \mathcal{L}_2$ und $\mathcal{L}_1 \subseteq \mathcal{L}_0$.

7. Rechtslineare Grammatiken und reguläre Sprachen

Satz 7.1

Die Typ-3-Sprachen sind genau die regulären/erkennbaren Sprachen, d.h.

$$\mathcal{L}_3 = \{L \mid L \text{ ist regulär}\}.$$

Beweis. Der Beweis wird in zwei Richtungen durchgeführt:

1. Jede Typ-3-Sprache ist erkennbar

Es sei $L \in \mathcal{L}_3$, d.h. $L = L(G)$ für eine Typ-3-Grammatik $G = (N, \Sigma, P, S)$. Es gilt $w \in L(G)$ gdw.

Es gibt eine Ableitung

$$(\star) \quad S = B_0 \vdash_G w_1 B_1 \vdash_G w_1 w_2 B_2 \vdash_G \dots \vdash_G w_1 \dots w_{n-1} B_{n-1} \vdash_G w_1 \dots w_{n-1} w_n$$

für Produktionen $B_{i-1} \rightarrow w_i B_i \in P$ ($i = 1, \dots, n$) und $B_{n-1} \rightarrow w_n \in P$.

Wir konstruieren nun einen NEA mit Worttransitionen, der die Nichtterminalsymbole von G als Zustände hat:

$\mathcal{A} = (N \cup \{\Omega\}, \Sigma, S, \Delta, \{\Omega\})$, wobei

- $\Omega \notin N$ neuer Endzustand ist und
- $\Delta = \{(A, w, B) \mid A \rightarrow wB \in P\} \cup \{(A, w, \Omega) \mid A \rightarrow w \in P\}$.

Ableitungen der Form (\star) entsprechen nun genau Pfaden in \mathcal{A} :

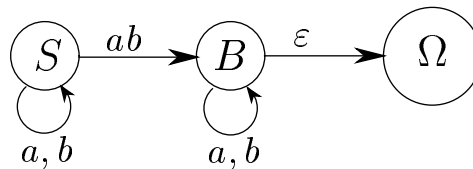
$$(\star\star) \quad (S, w_1, B_1)(B_1, w_2, B_2) \dots (B_{n-2}, w_{n-1}, B_{n-1})(B_{n-1}, w_n, \Omega)$$

Dies zeigt $L(\mathcal{A}) = L(G)$.

Beispiel 6.5 (Fortsetzung)

$$P = \{S \rightarrow aS, \\ S \rightarrow bS, \\ S \rightarrow abB, \\ B \rightarrow aB, \\ B \rightarrow bB, \\ B \rightarrow \varepsilon\}$$

liefert den folgenden NEA mit Wortübergängen:



2. Jede erkennbare Sprache ist eine Typ-3-Sprache

Es sei $L = L(\mathcal{A})$ für einen NEA $\mathcal{A} = (Q, \Sigma, q_0, \Delta, F)$. Wir definieren daraus eine Typ-3-Grammatik $G = (N, \Sigma, P, S)$ wie folgt:

$$\begin{aligned} N &:= Q \\ S &:= q_0 \\ P &:= \{p \longrightarrow aq \mid (p, a, q) \in \Delta\} \cup \{p \longrightarrow \varepsilon \mid p \in F\} \end{aligned}$$

Ein Pfad in \mathcal{A} der Form

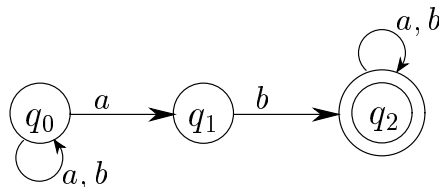
$$(q_0, a_1, q_1)(q_1, a_2, q_2) \dots (q_{n-1}, a_n, q_n)$$

mit $q_n \in F$ entspricht nun genau einer Ableitung

$$q_0 \vdash_G a_1 q_1 \vdash_G a_1 a_2 q_2 \vdash_G \dots \vdash_G a_1 \dots a_n q_n \vdash_G a_1 \dots a_n.$$

Beispiel:

Der folgende NEA



liefert die Grammatik mit den rechtslinearen Produktionen

$$\begin{aligned} P = \{ & q_0 \longrightarrow aq_0, \\ & q_0 \longrightarrow bq_0, \\ & q_0 \longrightarrow aq_1, \\ & q_1 \longrightarrow bq_2, \\ & q_2 \longrightarrow aq_2, \\ & q_2 \longrightarrow bq_2, \\ & q_2 \longrightarrow \varepsilon \} \end{aligned}$$

□

Korollar 7.2

$$\mathcal{L}_3 \subset \mathcal{L}_2.$$

Beweis. Wir wissen bereits, dass $\mathcal{L}_3 \subseteq \mathcal{L}_2$ gilt. Außerdem haben wir mit Beispiel 6.1 eine Sprache $L := \{a^n b^n \mid n \geq 0\} \in \mathcal{L}_2$. Im Teil I haben wir gezeigt, dass L nicht erkennbar/regulär ist, d.h. mit Satz 7.1 folgt $L \notin \mathcal{L}_3$. □

Beispiel 7.3

Als ein weiteres Beispiel für eine kontextfreie Sprache, die nicht regulär ist, betrachten wir $L = \{a^n b^m \mid n \neq m\}$. (Vgl. Beispiele 3.3, 4.2) Man kann diese Sprache mit folgender kontextfreier Grammatik erzeugen:

$G = (N, \Sigma, P, S)$ mit

- $N = \{S, S_{\geq}, S_{\leq}\}$
- $\Sigma = \{a, b\}$
- $P = \{S \rightarrow aS_{\geq}, \quad S \rightarrow S_{\leq}b,$
 $S_{\geq} \rightarrow aS_{\geq}b, \quad S_{\leq} \rightarrow aS_{\leq}b,$
 $S_{\geq} \rightarrow aS_{\geq}, \quad S_{\leq} \rightarrow S_{\leq}b,$
 $S_{\geq} \rightarrow \varepsilon, \quad S_{\leq} \rightarrow \varepsilon\}$

Es gilt nun:

- $S_{\geq} \vdash_G^* w \in \{a, b\}^* \Rightarrow w = a^n b^m$ mit $n \geq m$,
- $S_{\leq} \vdash_G^* w \in \{a, b\}^* \Rightarrow w = a^n b^m$ mit $n \leq m$,

woraus sich ergibt:

- $S \vdash_G^* w \in \{a, b\}^* \Rightarrow w = aa^n b^m$ mit $n \geq m$ oder $w = a^n b^m b$ mit $n \leq m$,

d.h. $L(G) = \{a^n b^m \mid n \neq m\}$.

8. Normalformen kontextfreier Grammatiken

Wir werden zeigen, dass man die Klasse aller kontextfreien Sprachen bereits mit kontextfreien Grammatiken einer *eingeschränkten* Form erzeugen kann.

Zwei Grammatiken heißen *äquivalent*, falls sie dieselbe Sprache erzeugen.

Zunächst zeigen wir, dass man „überflüssige“ Symbole aus kontextfreien Grammatiken eliminieren kann.

Definition 8.1 (terminierende, erreichbare Symbole; reduzierte Grammatik)

Es sei $G = (N, \Sigma, P, S)$ eine kontextfreie Grammatik.

- 1) $A \in N$ heißt *terminierend*, falls es ein $w \in \Sigma^*$ gibt mit $A \vdash_G^* w$.
- 2) $A \in N$ heißt *erreichbar*, falls es $u, v \in (\Sigma \cup N)^*$ gibt mit $S \vdash_G^* uAv$.
- 3) G heißt *reduziert*, falls alle Elemente von N *erreichbar* und *terminierend* sind.

Lemma 8.2

Zu einer kontextfreien Grammatik $G = (N, \Sigma, P, S)$ kann man effektiv die Menge der erreichbaren Nichtterminalsymbole bestimmen.

Beweis. Wir definieren dazu

$$E_0 := \{S\}$$

$$E_{i+1} := E_i \cup \{A \mid \exists B \in E_i \text{ mit Regel } B \longrightarrow u_1Au_2 \in P\}$$

Es gilt

$$E_0 \subseteq E_1 \subseteq E_2 \subseteq \dots \subseteq N.$$

Da N endlich ist, gibt es ein k mit $E_k = E_{k+1}$ und damit $E_k = \bigcup_{i \geq 0} E_i$.

Behauptung:

$E_k = \{A \in N \mid A \text{ ist erreichbar}\}$, denn:

„ \subseteq “: Offenbar enthalten alle E_i nur erreichbare Symbole

„ \supseteq “: Zeige durch Induktion über i : $S \vdash_G^i uAv \Rightarrow A \in E_i$.

□

Beispiel:

$$P = \left\{ \begin{array}{ll} S \longrightarrow aS, & A \longrightarrow ASB, \\ S \longrightarrow SB, & A \longrightarrow C, \\ S \longrightarrow SS, & B \longrightarrow Cb, \\ S \longrightarrow \varepsilon & \end{array} \right\}$$

$$E_0 = \{S\} \subset E_1 = \{S, B\} \subset E_2 = \{S, B, C\} = E_3$$

Es ist also A das einzige nichterreichbare Symbol.

Lemma 8.3

Zu einer kontextfreien Grammatik $G = (N, \Sigma, P, S)$ kann man effektiv die Menge der terminierenden Symbole bestimmen.

Beweis. Wir definieren dazu

$$T_1 := \{A \in N \mid \exists w \in \Sigma^* : A \longrightarrow w \in P\}$$

$$T_{i+1} := T_i \cup \{A \in N \mid \exists w \in (\Sigma \cup T_i)^* : A \longrightarrow w \in P\}$$

Es gilt

$$T_1 \subseteq T_2 \subseteq \dots \subseteq N.$$

Da N endlich ist, gibt es ein k mit $T_k = T_{k+1} = \bigcup_{i \geq 1} T_i$.

Behauptung:

$T_k = \{A \in N \mid A \text{ ist terminierend}\}$, denn:

„ \subseteq “: Offenbar sind alle Elemente von T_i ($i \geq 1$) terminierend:

- $i = 1$: $A \vdash_G w \in \Sigma^*$.
- $i \rightarrow i + 1$: $A \vdash_G u_1 B_1 u_2 B_2 \dots u_n B_n u_{n+1}$, $B_i \vdash_G^* w_i \in \Sigma^*$

„ \supseteq “: Zeige durch Induktion über i :

$$A \vdash_G^{\leq i} w \in \Sigma^* \Rightarrow A \in T_i.$$

- $i = 1$: $A \vdash_G^{\leq 1} w \in \Sigma^* \Rightarrow A \longrightarrow w \in P \Rightarrow A \in T_1$
- $i \rightarrow i + 1$: $A \vdash_G u_1 B_1 \dots u_n B_n u_{n+1}$
 $\vdash_G^{\leq i} w = u_1 w_1 \dots u_n w_n u_{n+1}$, $u_j w_j \in \Sigma^*$, $B_j \in N$
 \Rightarrow Für alle j gilt: $B_j \vdash_G^{\leq i} w_j \in \Sigma^*$
 $\Rightarrow B_j \in T_i$ (Induktion) $\Rightarrow A \in T_{i+1}$

□

Aus Lemma 8.3 folgt unmittelbar:

Satz 8.4

Das Leerheitsproblem ist für kontextfreie Sprachen entscheidbar.

Beweis. Offenbar gilt $L(G) \neq \emptyset$ gdw. $\exists w \in \Sigma^* : S \vdash_G^* w$ gdw. S ist terminierend. □

Lemma 8.2 und 8.3 zusammen zeigen, wie man unerreichbare und nichtterminierende Symbole eliminieren kann.

Satz 8.5

Zu jeder kontextfreien Grammatik G mit $L(G) \neq \emptyset$ kann man effektiv eine äquivalente reduzierte kontextfreie Grammatik erzeugen.

Beweis.

Erster Schritt: Eliminieren nicht terminierender Symbole.

Zu $G = (N, \Sigma, P, S)$ definieren wir $G' := (N', \Sigma, P', S)$, wobei

- $N' := \{A \in N \mid A \text{ ist terminierend in } G\}$
- $P' := \{A \longrightarrow w \in P \mid A \in N', w \in (N' \cup \Sigma)^*\}$

Beachte: Aus $L(G) \neq \emptyset$ folgt $S \in N'$!

Zweiter Schritt: Eliminieren unerreichbarer Symbole.

Wir definieren $G'' := (N'', \Sigma, P'', S)$, wobei

- $N'' := \{A \in N' \mid A \text{ ist erreichbar in } G'\}$
- $P'' := \{A \longrightarrow w \in P' \mid A \in N''\}$

Beachte: Ist $A \in N''$ und $A \longrightarrow u_1 B u_2 \in P'$, so ist $B \in N''$.

Man sieht leicht, dass $L(G) = L(G') = L(G'')$ und G'' reduziert ist. □

Vorsicht:

Die Reihenfolge der beiden Schritte ist wichtig! Symbole, die in G noch erreichbar waren, müssen es in G' nicht mehr sein.

Z.B.: $S \longrightarrow AB$ mit A terminierend, B nicht.

Als nächstes eliminieren wir Regeln der Form $A \longrightarrow \varepsilon$.

Lemma 8.6

Es sei G eine kontextfreie Grammatik. Dann lässt sich eine Grammatik G' ohne Regeln der Form $A \longrightarrow \varepsilon$ konstruieren mit $L(G') = L(G) \setminus \{\varepsilon\}$.

Beweis.

- 1) Finde alle $A \in N$ mit $A \vdash_G^* \varepsilon$:

$$N_1 := \{A \in N \mid A \longrightarrow \varepsilon \in P\}$$

$$N_{i+1} := N_i \cup \{A \in N \mid A \longrightarrow B_1 \dots B_n \in P \text{ für } B_j \in N_i\}$$

Es gibt ein k mit $N_k = N_{k+1} = \bigcup_{i \geq 1} N_i$. Für dieses k gilt: $A \in N_k \quad \text{gdw.} \quad A \vdash_G^* \varepsilon$.

- 2) Eliminiere in G alle Regeln $A \longrightarrow \varepsilon$. Um dies auszugleichen, nimmt man für alle Regeln

$$A \longrightarrow u_1 B_1 \dots u_n B_n u_{n+1} \text{ mit } B_i \in N_k \text{ und } u_i \in (\Sigma \cup N \setminus N_k)^*$$

die Regeln

$$A \longrightarrow u_1 \beta_1 u_2 \dots u_n \beta_n u_{n+1} \text{ mit } \beta_i \in \{B_i, \varepsilon\} \text{ und } u_1 \beta_1 u_2 \dots u_n \beta_n u_{n+1} \neq \varepsilon$$

hinzu. □

Beispiel:

$$\begin{aligned}
 P &= \{S \rightarrow aS, S \rightarrow SS, S \rightarrow bA, \\
 &\quad A \rightarrow BB, \\
 &\quad B \rightarrow CC, B \rightarrow aAbC, \\
 &\quad C \rightarrow \varepsilon\} \\
 N_0 &= \{C\}, \\
 N_1 &= \{C, B\}, \\
 N_2 &= \{C, B, A\} = N_3 \\
 P' &= \{S \rightarrow aS, S \rightarrow SS, S \rightarrow bA, S \rightarrow b, \\
 &\quad A \rightarrow BB, A \rightarrow B, \\
 &\quad B \rightarrow CC, B \rightarrow C, \\
 &\quad B \rightarrow aAbC, B \rightarrow abC, B \rightarrow aAb, B \rightarrow ab\}
 \end{aligned}$$

Die Ableitung $S \vdash bA \vdash bBB \vdash bCCB \vdash bCCCC \vdash^* b$ kann in G' direkt durch $S \vdash b$ erreicht werden.

Definition 8.7 (ε -freie kontextfreie Grammatik)

Eine kontextfreie Grammatik heißt ε -frei, falls gilt:

- 1) Sie enthält keine Regeln $A \rightarrow \varepsilon$ für $A \neq S$.
- 2) Ist $S \rightarrow \varepsilon$ enthalten, so kommt S nicht auf der rechten Seite einer Regel vor.

Satz 8.8

Zu jeder kontextfreien Grammatik G kann effektiv eine äquivalente ε -freie Grammatik konstruiert werden.

Beweis. Konstruiere G' wie im Beweis von Lemma 8.6 beschrieben. Ist $\varepsilon \notin L(G)$ (d.h. $S \not\vdash_G^* \varepsilon$, also $S \notin N_k$), so ist G' die gesuchte ε -freie Grammatik. Sonst erweitere G' um ein neues Startsymbol S_0 und die Produktionen $S_0 \rightarrow S$ und $S_0 \rightarrow \varepsilon$. □

Korollar 8.9

$$\mathcal{L}_2 \subseteq \mathcal{L}_1.$$

Beweis. Offenbar ist jede ε -freie kontextfreie Grammatik eine Typ-1-Grammatik:

- Produktionen: $u_1Au_2 \rightarrow u_1wu_2$ mit $|w| \geq 1$
- kontextfrei: $u_i = \varepsilon$; $S \rightarrow \varepsilon$ und S nicht auf rechter Seite). □

Wir zeigen nun, dass man auch auf Regeln der Form $A \rightarrow B$ verzichten kann.

Satz 8.10

Zu jeder kontextfreien Grammatik kann man effektiv eine äquivalente kontextfreie Grammatik konstruieren, die keine Regeln der Form $A \rightarrow B$ ($A, B \in N$) enthält.

Beweisskizze.

- 1) Bestimme zu jedem $A \in N$ die Menge $N(A) := \{B \in N \mid A \vdash_G^* B\}$
(effektiv machbar).
- 2) $P' = \{A \rightarrow w \mid B \rightarrow w \in P, B \in N(A) \text{ und } w \notin N\}$ □

Beispiel:

$$\begin{aligned}
 P &= \{S \rightarrow A, A \rightarrow B, B \rightarrow aA, B \rightarrow b\} \\
 N(S) &= \{S, A, B\}, \\
 N(A) &= \{A, B\}, \\
 N(B) &= \{B\} \\
 P' &= \{B \rightarrow aA, A \rightarrow aA, S \rightarrow aA, \\
 &\quad B \rightarrow b, A \rightarrow b, S \rightarrow b\}
 \end{aligned}$$

Wir zeigen nun, dass man sich auf Regeln sehr einfacher Form beschränken kann, nämlich $A \rightarrow a$ und $A \rightarrow BC$.

Satz 8.11 (Chomsky-Normalform)

Jede kontextfreie Grammatik lässt sich umformen in eine äquivalente Grammatik, die nur Regeln der Form

- $A \rightarrow a, A \rightarrow BC$ mit $A, B, C \in N, a \in \Sigma$
- und eventuell $S \rightarrow \varepsilon$, wobei S nicht rechts vorkommt

enthält. Eine derartige Grammatik heißt dann Grammatik in Chomsky-Normalform.

Beweis.

- 1) Konstruiere zu der gegebenen Grammatik eine äquivalente ε -freie ohne Regeln der Form $A \rightarrow B$. (Dabei ist die Reihenfolge wichtig!)
- 2) Führe für jedes $a \in \Sigma$ ein neues Nichtterminalsymbol X_a und die Produktion $X_a \rightarrow a$ ein.
- 3) Ersetze in jeder Produktion $A \rightarrow w$ mit $w \notin \Sigma$ alle Terminalsymbole a durch die zugehörigen X_a .
- 4) Produktionen $A \rightarrow B_1 \dots B_n$ für $n > 2$ werden ersetzt durch

$$A \rightarrow B_1 C_1, C_1 \rightarrow B_2 C_2, \dots, C_{n-2} \rightarrow B_{n-1} B_n$$

wobei die C_i jeweils neue Symbole sind.

□

Beachte:

Für einen NEA \mathcal{A} (eine rechtslineare Grammatik G) kann man das Wortproblem, also die Frage „ $w \in L(\mathcal{A})?$ “ („ $w \in L(G)?$ “) z.B. deshalb entscheiden, da ein Pfad in \mathcal{A} , der w erkennt (eine Ableitung in G , die w erzeugt) *genau* die Länge $|w|$ haben muss. Es gibt aber nur endlich viele Pfade (Ableitungen) dieser festen Länge.

Bei allgemeinen kontextfreien Grammatiken kann man keine Schranke für die Länge einer Ableitung von w aus S angeben, wohl aber bei Grammatiken in Chomsky-Normalform:

- Produktionen der Form $A \rightarrow BC$ verlängern um 1, d.h. sie können maximal $|w| - 1$ -mal angewandt werden.
- Produktionen der Form $A \rightarrow a$ erzeugen genau ein Terminalsymbol von w , d.h. sie werden genau $|w|$ -mal angewandt.

Es folgt: $w \in L(G) \setminus \{\varepsilon\}$ wird durch eine Ableitung der Länge $2|w| - 1$ erzeugt.

Da es aber i.a. $\geq 2^n$ Ableitungen der Länge n geben kann, liefert dies ein *exponentielles Verfahren* zur Entscheidung des Wortproblems.

Ein besseres Verfahren (*kubisch*) liefert die folgende Überlegung:

Definition 8.12

Es sei $G = (N, \Sigma, P, S)$ eine kontextfreie Grammatik in Chomsky-Normalform und $w = a_1 \dots a_n \in \Sigma^*$. Wir definieren:

- $w_{ij} := a_i \dots a_j$ (für $i \leq j$)
- $N_{ij} := \{A \in N \mid A \vdash_G^* w_{ij}\}$

Mit dieser Notation gilt nun:

- | | | |
|----|----------------------------|---|
| 1) | $S \in N_{1n}$ | <u>gdw.</u> $w \in L(G)$ |
| 2) | $A \in N_{ii}$ | <u>gdw.</u> $A \vdash_G^* a_i$ |
| | | <u>gdw.</u> $A \rightarrow a_i \in P$ |
| 3) | $A \in N_{ij}$ für $i < j$ | <u>gdw.</u> $A \vdash_G^* a_i \dots a_j$ |
| | | <u>gdw.</u> $\exists A \rightarrow BC \in P$ und ein k mit $i \leq k < j$ mit |
| | | $B \vdash_G^* a_i \dots a_k$ und $C \vdash_G^* a_{k+1} \dots a_j$ |
| | | <u>gdw.</u> $\exists A \rightarrow BC \in P$ und k mit $i \leq k < j$ mit |
| | | $B \in N_{ik}$ und $C \in N_{(k+1)j}$ |

Diese Überlegungen liefern das folgende iterative Verfahren zur Bestimmung von N_{1k} :

Algorithmus 8.13 (CYK-Algorithmus von Cocke, Younger, Kasami)

```

FOR  $i := 1$  TO  $n$  DO
     $N_{ii} := \{A \mid A \rightarrow a_i \in P\}$ 
FOR  $d := 1$  TO  $n - 1$  DO    (wachsende Differenz  $d = j - i$ )
    FOR  $i := 1$  TO  $n - d$  DO
         $j := i + d$ 
         $N_{ij} := \emptyset$ 
        FOR  $k := i$  TO  $j - 1$  DO
             $N_{ij} := N_{ij} \cup \{A \mid \exists A \rightarrow BC \in P \text{ mit } B \in N_{ik} \text{ und } C \in N_{(k+1)j}\}$ 
    
```

Beachte:

In der innersten Schleife sind die Differenzen $k - i$ und $j - k + 1$ kleiner als das aktuelle d , also sind N_{ik} und $N_{(k+1)j}$ bereits berechnet.

Satz 8.14

Für eine gegebene Grammatik in Chomsky-Normalform entscheidet Algorithmus 8.13 die Frage „ $w \in L(G)$?“ in der Zeit $O(|w|^3)$.

Beweis. Drei geschachtelte Schleifen, die jeweils $\leq |w| = n$ Schritte machen, daraus folgt: $|w|^3$ Schritte in der innersten Schleife. □

Beachte:

Die Größe von G ist hier als konstant angenommen (fest vorgegebenes G). Daher ist die Suche nach den Produktionen $A \rightarrow BC$ und $A \rightarrow a_i$ auch konstant.

Beispiel:

$P = \{S \rightarrow SA, S \rightarrow a,$
 $A \rightarrow BS,$
 $B \rightarrow BB, B \rightarrow BS, B \rightarrow b, B \rightarrow c\}$

und $w = abacba$:

$i \backslash j$	1	2	3	4	5	6
1	S	\emptyset	S	\emptyset	\emptyset	S
2		B	A, B	B	B	A, B
3			S	\emptyset	\emptyset	S
4				B	B	A, B
5					B	A, B
6						S
$w =$	a	b	a	c	b	a

$S \in N_{1,6} = \{S\}$
 $\Rightarrow w \in L(G)$

Eine weitere interessante Normalform für kontextfreie Grammatiken ist die *Greibach-Normalform*, bei der jede Produktion mindestens ein Terminalsymbol erzeugt.

Satz 8.15

Jede kontextfreie Grammatik lässt sich effektiv umformen in eine äquivalente Grammatik, die nur Regeln der Form

- $A \longrightarrow aw$ ($A \in N, a \in \Sigma, w \in N^*$)
- und eventuell $S \longrightarrow \varepsilon$, wobei S nicht rechts vorkommt

enthält (ohne Beweis!).

Eine derartige Grammatik heißt Grammatik in Greibach-Normalform.

9. Abschlusseigenschaften kontextfreier Sprachen

Satz 9.1

Die Klasse \mathcal{L}_2 der kontextfreien Sprachen ist unter Vereinigung, Konkatenation und Kleene-Stern abgeschlossen.

Beweis. Es seien $L_1 = L(G_1)$ und $L_2 = L(G_2)$ die Sprachen für kontextfreie Grammatiken $G_i = (N_i, \Sigma, P_i, S_i)$ ($i = 1, 2$). O.B.d.A. nehmen wir an, dass $N_1 \cap N_2 = \emptyset$.

- 1) $G := (N_1 \cup N_2 \cup \{S\}, \Sigma, P_1 \cup P_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\}, S)$ mit $S \notin N_1 \cup N_2$ ist eine kontextfreie Grammatik mit $L(G) = L_1 \cup L_2$.
- 2) $G' := (N_1 \cup N_2 \cup \{S\}, \Sigma, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\}, S)$ mit $S \notin N_1 \cup N_2$ ist eine kontextfreie Grammatik mit $L(G') = L_1 \cdot L_2$.
- 3) $G'' := (N_1 \cup \{S\}, \Sigma, P_1 \cup \{S \rightarrow \varepsilon, S \rightarrow S S_1\}, S)$ mit $S \notin N_1$ ist eine kontextfreie Grammatik für L_1^* □

Wir werden zeigen, dass Abschluss unter *Durchschnitt* und *Komplement nicht* gilt. Dazu benötigen wir zunächst eine geeignete Methode, von einer Sprache nachzuweisen, dass sie *nicht* kontextfrei ist. Dies gelingt wieder mit Hilfe eines Pumping-Lemmas. Um dieses zu zeigen, stellt man Ableitungen als Bäume, sogenannte *Ableitungsbäume* dar.

Beispiel:

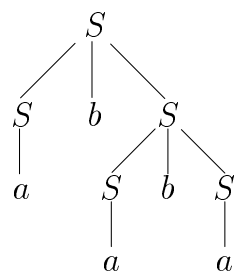
$$P = \{S \rightarrow SbS, S \rightarrow a\}$$

Drei *Ableitungen* des Wortes *ababa*:

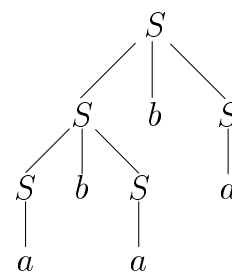
- 1) $S \vdash SbS \vdash abS \vdash abSbS \vdash ababS \vdash ababa$
- 2) $S \vdash SbS \vdash abS \vdash abSbS \vdash abSba \vdash ababa$
- 3) $S \vdash SbS \vdash Sba \vdash SbSba \vdash Sbaba \vdash ababa$

Die zugehörigen *Ableitungsbäume*:

Für 1) und 2):



Für 3):



Ein Ableitungsbaum kann also für mehr als eine Ableitung stehen und dasselbe Wort kann verschiedene Ableitungsbäume haben.

Allgemein:

Die Knoten des Ableitungsbaumes sind mit Elementen aus $\Sigma \cup N$ beschriftet. Ein mit A beschrifteter Knoten kann mit $\alpha_1, \dots, \alpha_n$ beschriftete Nachfolgerknoten haben, wenn $A \rightarrow \alpha_1 \dots \alpha_n \in P$ ist.

Ein Ableitungsbaum, dessen Wurzel mit A beschriftet ist und dessen Blätter (von links nach rechts) mit $\alpha_1, \dots, \alpha_n \in \Sigma \cup N$ beschriftet sind, beschreibt eine Ableitung $A \vdash_G^* \alpha_1 \dots \alpha_n$.

Lemma 9.2 (Pumping-Lemma für kontextfreie Sprachen)

Es sei G eine ε -freie kontextfreie Grammatik, die

- keine Regeln der Form $A \rightarrow B$ enthält,
- m nichtterminale Symbole enthält
- nur rechten Regelseiten der Länge $\leq k$ hat

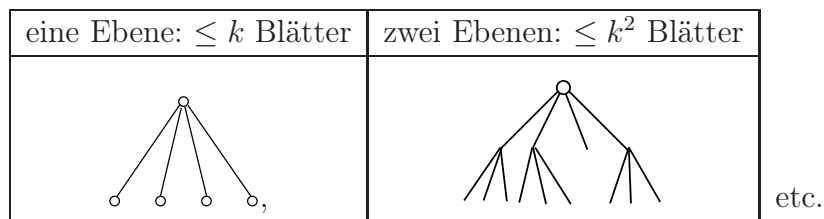
Es sei $n = k^{m+1}$.

Dann gibt es für jedes $z \in L(G)$ mit $|z| > n$ eine Zerlegung $z = uvwxy$ mit:

- $vx \neq \varepsilon$ und $|vwx| \leq n$
- $uw^iwx^iy \in L(G)$ für alle $i \geq 0$.

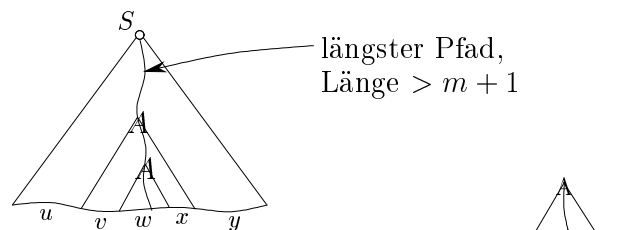
Beweis.

- 1) Ein Baum der Tiefe $\leq t$ und der Verzweigungszahl $\leq k$ hat maximal k^t viele Blätter:



Da der Ableitungsbaum für z als Blattanzahl $|z| > k^{m+1}$ hat, ist die maximale Tiefe (längster Pfad von Wurzel bis Blatt) $> m + 1$.

- 2) Da es nur m verschiedene Elemente in N gibt, kommt auf diesem längsten Pfad ein nichtterminales Symbol A zweimal vor.



Wir wählen hier o.B.d.A. A so, dass es in dem Baum keine andere Variablenwiederholung gibt. Daher hat dieser Baum eine Tiefe $\leq m + 1$, was $|vwx| \leq k^{m+1} = n$ zeigt.

3) Es gilt: $S \vdash_G^* uAy$, $A \vdash_G^* vAx$, $A \vdash_G^* w$, woraus folgt:

$$S \vdash_G^* uAy \vdash_G^* uv^iAx^i y \vdash_G^* uv^iwx^i y.$$

4) $vx \neq \varepsilon$: Da G ε -frei ist, wäre sonst $A \vdash_G^* vAx$ nur bei Anwesenheit von Regeln der Form $A \rightarrow B$ möglich. \square

Satz 9.3

$\mathcal{L}_2 \subset \mathcal{L}_1$.

Beweis. Wir haben bereits gezeigt, dass $\mathcal{L}_2 \subseteq \mathcal{L}_1$ gilt (Korollar 8.9). Es bleibt zu zeigen, dass die Inklusion echt ist.

Dafür betrachten wir die Sprache $L = \{a^n b^n c^n \mid n \geq 1\}$, und zeigen: $L \in \mathcal{L}_1 \setminus \mathcal{L}_2$

1) Wir zeigen zunächst $L \notin \mathcal{L}_2$. Angenommen, $L \in \mathcal{L}_2$. Dann gibt es eine ε -freie kontextfreie Grammatik G ohne Regeln der Form $A \rightarrow B$ für L . Es sei $n_0 = k^{m+1}$ die zugehörige Zahl aus Lemma 9.2. Wir betrachten $z = a^{n_0} b^{n_0} c^{n_0} \in L = L(G)$. Mit Satz 9.2 gibt es eine Zerlegung

$$z = uvwxy, \quad vx \neq \varepsilon \text{ und } uv^iwx^i y \in L \text{ für alle } i \geq 0.$$

1. Fall: v enthält verschiedene Symbole. Man sieht leicht, dass dann

$$uv^2wx^2y \notin a^*b^*c^* \supseteq L.$$

2. Fall: x enthält verschiedene Symbole. Dies führt zu entsprechendem Widerspruch.

3. Fall: v enthält lauter gleiche Symbole und x enthält lauter gleiche Symbole. Dann gibt es einen Symbole aus $\{a, b, c\}$, der in xv nicht vorkommt. Daher kommt dieser in $uv^0wx^0y = uwy$ weiterhin n_0 -mal vor. Aber es gilt $|uwy| < 3n_0$, was $uwy \notin L$ zeigt.

2) In Beispiel 6.3 haben wir eine Grammatik G mit $L(G) = \{a^n b^n c^n \mid n \geq 1\}$ angegeben. Leider enthält G eine Produktion, die nicht die Bedingungen für kontextsensitive Produktionen ($uAv \rightarrow uvv$, $|w| \geq 1$) erfüllt: $cB \rightarrow Bc$. Wir modifizieren G wie folgt:

- Ersetze in allen Produktionen c durch ein neues nichtterminales Symbol C und nimm die Produktion $C \rightarrow c$ hinzu:

$$\{S \rightarrow aSBC, S \rightarrow abC, CB \rightarrow BC, bB \rightarrow bb, C \rightarrow c\}.$$

- Ersetze nun $CB \rightarrow BC$ durch die kontextsensitiven Produktionen

$$CB \rightarrow A_1B, A_1B \rightarrow A_1A_2, A_1A_2 \rightarrow BA_2, BA_2 \rightarrow BC.$$

Diese Produktionen können nur dazu verwendet werden, $CB \rightarrow BC$ zu simulieren. \square

Beachte:

Auf diese Art kann man leicht zeigen, dass *jede nichtkürzende* Produktion $u \rightarrow v$ (mit $|u| \leq |v|$) durch *kontextsensitive* Produktionen ersetzt werden kann, ohne die Sprache zu ändern.

Korollar 9.4

Die Klasse \mathcal{L}_2 der kontextfreien Sprachen ist nicht unter Durchschnitt und Komplement abgeschlossen.

Beweis.

1) Die Sprachen $\{a^n b^n c^m \mid n \geq 1, m \geq 1\}$ und $\{a^m b^n c^n \mid n \geq 1, m \geq 1\}$ sind in \mathcal{L}_2 :

$$\bullet \{a^n b^n c^m \mid n \geq 1, m \geq 1\} = \underbrace{\{a^n b^n \mid n \geq 1\}}_{\in \mathcal{L}_2} \cdot \underbrace{\{c^m \mid m \geq 1\}}_{= c^+ \in \mathcal{L}_3 \subseteq \mathcal{L}_2}$$

$$\underbrace{\hspace{15em}}_{\in \mathcal{L}_2 \text{ (Konkatenation)}}$$

• $\{a^m b^n c^n \mid n \geq 1, m \geq 1\}$ — analog

2) $\{a^n b^n c^n \mid n \geq 1\} = \{a^n b^n c^m \mid n, m \geq 1\} \cap \{a^m b^n c^n \mid n, m \geq 1\}$.

Wäre \mathcal{L}_2 unter \cap abgeschlossen, so würde $\{a^n b^n c^n \mid n \geq 1\} \in \mathcal{L}_2$ folgen.

Widerspruch zu Teil 1) des Beweises von Satz 9.3.

3) $L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$.

Wäre \mathcal{L}_2 unter Komplement abgeschlossen, so auch unter \cap , da \mathcal{L}_2 unter \cup abgeschlossen ist. Widerspruch zu 2).

□

Beachte:

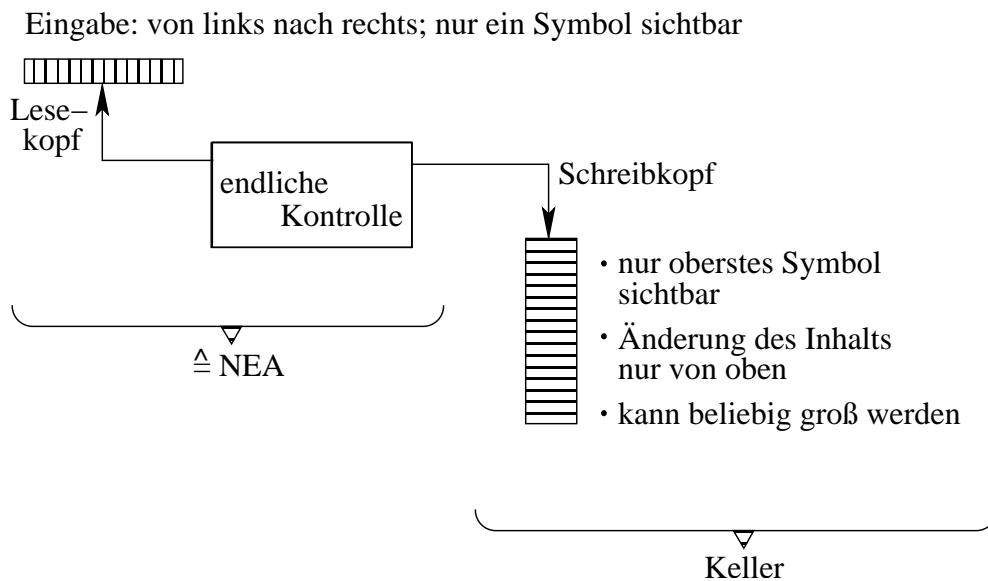
Man kann daher das Äquivalenzproblem für kontextfreie Sprachen nicht einfach auf das Leerheitsproblem reduzieren. Wir werden später sehen, dass das Äquivalenzproblem für kontextfreie Sprachen sogar *unentscheidbar* ist.

10. Kellerautomaten

Bisher hatten wir kontextfreie Sprachen nur mit Hilfe von Grammatiken charakterisiert. Wir haben gesehen, dass endliche Automaten *nicht* in der Lage sind, alle kontextfreien Sprachen zu akzeptieren.

Um die *Beschreibung von kontextfreien Sprachen* mit Hilfe von endlichen Automaten zu ermöglichen, muss man diese um eine (potentiell unendliche) *Speicherkomponente*, einen sogenannten *Keller*, erweitern.

Die folgende Abbildung zeigt eine schematische Darstellung eines *Kellerautomaten*:



Definition 10.1 (Kellerautomat)

Ein *Kellerautomat* (*pushdown automaton*, kurz *PDA*) hat die Form

$\mathcal{A} = (Q, \Sigma, \Gamma, q_0, Z_0, \Delta)$, wobei

- Q eine endliche Menge von *Zuständen* ist,
- Σ das *Eingabealphabet* ist,
- Γ das *Kelleralphabet* ist,
- $q_0 \in Q$ der *Anfangszustand* ist,
- $Z_0 \in \Gamma$ das *Kellerstartsymbol* ist und
- $\Delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma \times \Gamma^* \times Q$ eine endliche *Übergangsrelation* ist.

Anschaulich bedeutet die Übergangsrelation:

(q, a, Z, γ, q') : Im Zustand q mit aktuellem Eingabesymbol a und oberstem Kellersymbol Z darf der Automat Z durch γ ersetzen und in den Zustand q' und zum nächsten Eingabesymbol übergehen.

$(q, \varepsilon, Z, \gamma, q')$: wie oben, nur dass das aktuelle Eingabesymbol nicht relevant ist und man nicht zum nächsten Eingabesymbol übergeht (der Lesekopf ändert seine Position nicht).

Den aktuellen *Zustand* einer Kellerautomatenberechnung kann man beschreiben durch

- den *noch zu lesenden Rest* $w \in \Sigma^*$ der Eingabe (Lesekopf steht auf dem ersten Symbol von w)
- den *Zustand* $q \in Q$
- den *Kellerinhalt* $\gamma \in \Gamma^*$ (Schreiblesekopf steht auf dem ersten Symbol von γ)

Definition 10.2

Eine *Konfiguration* von \mathcal{A} hat die Form

$$\mathcal{K} = (q, w, \gamma) \in Q \times \Sigma^* \times \Gamma^*.$$

Die Übergangsrelation ermöglicht die folgenden *Konfigurationsübergänge*:

- $(q, aw, Z\gamma) \vdash (q', w, \beta\gamma)$ falls $(q, a, Z, \beta, q') \in \Delta$
- $(q, w, Z\gamma) \vdash (q', w, \beta\gamma)$ falls $(q, \varepsilon, Z, \beta, q') \in \Delta$
- $\mathcal{K} \vdash^* \mathcal{K}'$ gdw. $\exists n \geq 0 \exists \mathcal{K}_0, \dots, \mathcal{K}_n$ mit

$$\mathcal{K}_0 = \mathcal{K}, \mathcal{K}_n = \mathcal{K}' \text{ und } \mathcal{K}_i \vdash \mathcal{K}_{i+1} \text{ für } 0 \leq i < n.$$

Der Automat \mathcal{A} *akzeptiert* das Wort $w \in \Sigma^*$ gdw.

$$(q_0, w, Z_0) \vdash^* (q, \varepsilon, \varepsilon) \quad (\text{Eingabewort ganz gelesen und Keller leer}).$$

Die von \mathcal{A} *akzeptierte Sprache* ist

$$L(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ akzeptiert } w\}.$$

Beispiel 10.3

Ein PDA für $\{a^n b^n \mid n \geq 1\}$.

- $Q = \{q_0, q_1, f\}$,
- $\Gamma = \{Z, Z_0\}$,
- $\Sigma = \{a, b\}$ und
- $\Delta = \{(q_0, a, Z_0, ZZ_0, q_0), \text{ (erstes } a, \text{ speichere } Z)$
 $(q_0, a, Z, ZZ, q_0), \text{ (weitere } a\text{'s, speichere } Z)$
 $(q_0, b, Z, \varepsilon, q_1), \text{ (erstes } b, \text{ entnimm } Z)$
 $(q_1, b, Z, \varepsilon, q_1), \text{ (weitere } b\text{'s, entnimm } Z)$
 $(q_1, \varepsilon, Z_0, \varepsilon, f)\}$ (lösche das Kellerstartsymbol)

Betrachten wir uns einige Konfigurationsübergänge:

- 1) $(q_0, aabb, Z_0) \vdash (q_0, abb, ZZ_0) \vdash (q_0, bb, ZZZ_0) \vdash (q_1, b, ZZ_0) \vdash (q_1, \varepsilon, Z_0) \vdash (f, \varepsilon, \varepsilon)$
– akzeptiert
- 2) $(q_0, aab, Z_0) \vdash^* (q_0, b, ZZZ_0) \vdash (q_1, \varepsilon, ZZ_0)$
– kein Übergang mehr möglich
- 3) $(q_0, abb, Z_0) \vdash (q_0, bb, ZZ_0) \vdash (q_1, b, Z_0) \vdash (f, b, \varepsilon)$
– nicht akzeptiert

Der Kellerautomat aus Beispiel 10.3 ist deterministisch, da es zu jeder Konfiguration höchstens eine Folgekonfiguration gibt.

Definition 10.4

Der PDA $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, Z_0, \Delta)$ heißt *deterministisch*, falls die folgenden beiden Bedingungen erfüllt sind:

- 1) $\forall q \in Q \forall a \in \Sigma \cup \{\varepsilon\} \forall Z \in \Gamma$ existiert höchstens ein Übergang der Form $(q, a, Z, \dots, \dots) \in \Delta$.
- 2) Existiert ein Tupel $(q, \varepsilon, Z, \dots, \dots) \in \Delta$, so existiert kein Tupel der Form $(q, a, Z, \dots, \dots) \in \Delta$ mit $a \in \Sigma$.

Beispiel:

Die Sprache $L = \{w \bar{w} \mid w \in \{a, b\}^*\}$ (wobei für $w = a_1 \dots a_n$ gilt $\bar{w} = a_n \dots a_1$) wird von einem nichtdeterministischen PDA akzeptiert.

Idee: Der PDA legt die erste Worthälfte im Keller ab und greift darauf *in umgekehrter Reihenfolge* beim Lesen der zweiten Worthälfte zurück. So kann Nichtübereinstimmung festgestellt werden (kein Übergang bei Nichtübereinstimmung).

Nichtdeterminismus ist intuitiv nötig, da man „raten“ muss, wann die Wortmitte erreicht ist (Man kann beweisen, dass kein deterministischer PDA diese Sprache akzeptiert).

$$Q = \{q_0, q_1, q_2, f\}, \Gamma = \{a, b, Z_0\}, \Sigma = \{a, b\}, \Delta =$$

$\{(q_0, \varepsilon, Z_0, \varepsilon, f),$	akzeptiert ε
$(q_0, a/b, Z_0, a/bZ_0, q_1),$	Lesen
$(q_1, a, b, ab, q_1),$	und
$(q_1, b, a, ba, q_1),$	Speichern
$(q_1, a, a, aa, q_1),$	der ersten
$(q_1, b, b, bb, q_1),$	Worthälfte
$(q_1, b, b, \varepsilon, q_2),$	Lesen der
$(q_1, a, a, \varepsilon, q_2),$	zweiten Hälfte
$(q_2, b, b, \varepsilon, q_2),$	und Vergleichen
$(q_2, a, a, \varepsilon, q_2),$	mit erster
$(q_2, \varepsilon, Z_0, \varepsilon, f)\}$	Entfernen Kellerstartsymbol

Es gibt auch einen anderen Akzeptanzbegriff für PDAs, bei dem man nicht fordert, dass der Keller leer ist, sondern dass ein Endzustand (aus einer Menge $F \subseteq Q$) erreicht ist. Dieser Akzeptanzbegriff liefert aber dieselbe Sprachklasse.

Wir werden nun zeigen, dass man mit Kellerautomaten genau die kontextfreien Sprachen akzeptieren kann.

Satz 10.5

Für eine formale Sprache L sind äquivalent:

- 1) $L = L(G)$ für eine kontextfreie Grammatik G .
- 2) $L = L(\mathcal{A})$ für einen PDA \mathcal{A} .

Beweis.

„1 \rightarrow 2“: Es sei $G = (N, \Sigma, P, S)$ eine kontextfreie Grammatik. Der zugehörige PDA simuliert Linksableitungen von G , d.h. Ableitungen, bei denen stets das am weitesten links stehende Nichtterminalsymbol zuerst ersetzt wird.

Beachte:

Jedes Wort in $L(G)$ kann durch eine Linksableitung erzeugt werden (da G kontextfrei ist).

Wir definieren dazu $\mathcal{A} = (\{q\}, \Sigma, \Sigma \cup N, q, S, \Delta)$ mit $\Delta :=$

- $\{(q, \varepsilon, A, \gamma, q) \mid A \rightarrow \gamma \in P\} \cup$ (Anwenden einer Produktion auf oberstes Kellersymbol) (★)
- $\{(q, a, a, \varepsilon, q) \mid a \in \Sigma\}$ (Entfernen bereits erzeugter Terminalsymbole von der Kellerspitze, wenn sie in der Eingabe vorhanden sind) (★★)

Beispiel:

$P = \{S \rightarrow \varepsilon, S \rightarrow aSa, S \rightarrow bSb\}$ liefert die Übergänge:

- $(q, \varepsilon, S, \varepsilon, q), (S \rightarrow \varepsilon)$
- $(q, \varepsilon, S, aSa, q), (S \rightarrow aSa)$
- $(q, \varepsilon, S, bSb, q), (S \rightarrow bSb)$
- $(q, a, a, \varepsilon, q), (a \text{ entfernen})$
- $(q, b, b, \varepsilon, q), (b \text{ entfernen})$

Die Ableitung $S \vdash aSa \vdash abSba \vdash abba$ entspricht der *Konfigurationsfolge*

- $(q, abba, S) \vdash (q, abba, aSa) \vdash (q, bba, Sa) \vdash (q, bba, bSba) \vdash$
- $(q, ba, Sba) \vdash (q, ba, ba) \vdash (q, a, a) \vdash (q, \varepsilon, \varepsilon)$

Behauptung:

Für $u, v \in \Sigma^*$ und $\alpha \in \{\varepsilon\} \cup N \cdot (\Sigma \cup N)^*$ gilt:

$$S \vdash_G^* u\alpha \text{ mit Linksableitung} \quad \text{gdw.} \quad (q, uv, S) \vdash^* (q, v, \alpha).$$

Beachte:

Für $\alpha = \varepsilon = v$ folgt:

$$S \vdash_G^* u \text{ gdw. } (q, u, S) \vdash^* (q, \varepsilon, \varepsilon)$$

d.h. $L(G) = L(\mathcal{A})$.

Beweis der Behauptung.

„ \Leftarrow “: Beweis durch Induktion über die Anzahl k der Übergänge mit Transitionen der Form $(\star) : (q, \varepsilon, A, \gamma, q)$.

$k = 0$:

Da im Keller S steht, sind auch keine Transitionen der Form $(\star\star) : (q, a, a, \varepsilon, q)$ möglich, d.h. $u = \varepsilon$ und $\alpha = S$. Offenbar gilt $S \vdash^* S$.

$k \rightarrow k + 1$:

$$(q, \underbrace{u_1 u_2}_u v, S) \vdash^* (q, u_2 v, A\alpha') \xrightarrow{\text{letzte Trans. } (\star)} \vdash (q, u_2 v, \underbrace{\gamma\alpha'}_{u_2\alpha}) \xrightarrow{\text{Trans. } (\star\star)} \vdash^* (q, v, \alpha)$$

Induktion liefert: $S \vdash^* u_1 A\alpha'$ und außerdem wegen $A \rightarrow \gamma \in P$:

$$u_1 A\alpha' \vdash u_1 \gamma\alpha' = u_1 u_2 \alpha = u\alpha.$$

„ \Rightarrow “: Beweis durch Induktion über die Länge der Linksableitung

$k = 0$:

Dann ist $u = \varepsilon$, $\alpha = S$ und $(q, v, S) \vdash^0 (q, v, S)$

$k \rightarrow k + 1$:

$S \vdash_G^{k+1} u\alpha$, d.h.

$S \vdash_G^k u' A\beta \vdash \underbrace{u'\gamma\beta}_{=u\alpha}$ mit $u' \in \Sigma^*$ (da Linksableitung) und $A \rightarrow \gamma \in P$.

Es sei u'' das längste Anfangsstück von $\gamma\beta$, das in Σ^* liegt. Dann ist $u'u'' = u$ und α ist der Rest von $\gamma\beta$, d.h. $\gamma\beta = u''\alpha$.

Induktion liefert:

$$(q, u' \underbrace{u''v}_{v'}, S) \vdash^* (q, u''v, A\beta) \vdash (q, u''v, \gamma\beta) = (q, u''v, u''\alpha).$$

Außerdem gibt es Übergänge $(q, u''v, u''\alpha) \vdash^* (q, v, \alpha)$. □

„ $2 \rightarrow 1$ “: Es sei $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, Z_0, \Delta)$ ein PDA. Die Nichtterminalsymbole der zugehörigen Grammatik G sind alle Tripel $[p, Z, q] \in Q \times \Gamma \times Q$.

Idee:

Es soll gelten: $[p, Z, q] \vdash_G^* u \in \Sigma^*$ gdw.

- A kommt vom Zustand p in den Zustand q
- durch *Lesen von u* auf dem Eingabeband und
- *Löschen von Z* aus dem Keller

(ohne dabei die Symbole unter Z anzutasten).

Wie realisiert man dies durch geeignete Produktionen?

Betrachte den PDA-Übergang $(p, a, Z, X_1 \dots X_n, p') \in \Delta$. Hier wird a auf dem Eingabeband verbraucht (falls nicht $a = \varepsilon$). Z wird durch $X_1 \dots X_n$ ersetzt.

Um den Kellerinhalt zu erreichen, den man durch einfaches Wegstreichen von Z erhalten würde, muss man also nun noch $X_1 \dots X_n$ löschen. Löschen von X_i kann durch das Nichtterminalsymbol $[p_{i-1}, X_i, p_i]$ (für geeignete Zwischenzustände p_{i-1}, p_i) beschrieben werden.

Formale Definition:

$$\begin{aligned}
 G &:= (N, \Sigma, P, S) \text{ mit} \\
 N &:= \{S\} \cup \{[p, Z, q] \mid p, q \in Q, Z \in \Gamma\} \\
 P &:= \{S \longrightarrow [q_0, Z_0, q] \mid q \in Q\} \cup \\
 &\quad \{[p, Z, q] \longrightarrow a[p_0, X_1, p_1][p_1, X_2, p_2] \dots [p_{n-1}, X_n, q] \mid \\
 &\quad p, q, p_0, p_1, \dots, p_{n-1} \in Q, \\
 &\quad a \in \Sigma \cup \{\varepsilon\} \text{ und} \\
 &\quad (p, a, Z, X_1 \dots X_n, p_0) \in \Delta \text{ und} \\
 &\quad q = p_0 \text{ falls } n = 0\}
 \end{aligned}$$

Beachte:

Für $n = 0$ hat man die Produktion $[p, Z, q] \longrightarrow a$, welche dem Übergang $(p, a, Z, \varepsilon, q) \in \Delta$ entspricht.

Behauptung:

Für alle $p, q \in Q, u \in \Sigma^*, Z \in \Gamma$ gilt:

$$(*) \quad (p, u, Z) \vdash_{\mathcal{A}}^* (q, \varepsilon, \varepsilon) \text{ gdw. } [p, Z, q] \vdash_G^* u$$

Für $p = p_0$ und $Z = Z_0$ folgt daraus:

$$(q_0, u, Z_0) \vdash^* (q, \varepsilon, \varepsilon) \text{ gdw. } S \vdash_G [q_0, Z_0, q] \vdash_G^* u$$

d.h. $u \in L(\mathcal{A})$ gdw. $u \in L(G)$.

Der Beweis dieser Behauptung kann durch Induktion über die Länge der Konfigurationsfolge („ \Rightarrow “) bzw. über die Länge der Ableitung („ \Leftarrow “) geführt werden.

□

Beispiel: (vgl. Beispiel 10.3)

Gegeben ist der PDA für $\{a^n b^n \mid n \geq 1\}$.

$$(q_0, ab, Z_0) \xrightarrow{(q_0, a, Z_0, Z Z_0, q_0)} (q_0, b, Z Z_0) \xrightarrow{(q_0, b, Z, \varepsilon, q_1)} (q_1, \varepsilon, Z_0) \xrightarrow{(q_1, \varepsilon, Z_0, \varepsilon, f)} (f, \varepsilon, \varepsilon)$$

entspricht der Ableitung

$$S \vdash_G [q_0, Z_0, f] \vdash_G a[q_0, Z, q_1][q_1, Z_0, f] \vdash_G ab[q_1, Z_0, f] \vdash_G ab.$$

III. Berechenbarkeit

Einführung

Aus der Sicht der Theorie der *formalen Sprachen* geht es in diesem Teil darum, Automatenmodelle für die Typ-0- und die Typ-1-Sprachen zu finden. Allgemeiner geht es aber darum, die intuitiven Begriffe „*berechenbare Funktion*“ und „*entscheidbares Problem*“ zu formalisieren.

Um für eine (eventuell partielle) Funktion

$$f : \mathbb{N}^k \rightarrow \mathbb{N} \quad (\text{bzw. } (\Sigma^*)^k \rightarrow \Sigma^*)$$

intuitiv klarzumachen, dass sie berechenbar ist, genügt es, ein Berechnungsverfahren (einen *Algorithmus*) für die Funktion anzugeben (z.B. in Form eines Modula-, Pascal- oder C-Programmes oder einer *abstrakten Beschreibung* der Vorgehensweise bei der Berechnung). Entsprechend hatten wir in den Teilen I und II die Entscheidbarkeit von Problemen (Wortproblem, Leerheitsproblem, Äquivalenzproblem, ...) dadurch begründet, dass wir *intuitiv* beschrieben haben, wie man die Probleme mit Hilfe eines Rechenverfahrens (d.h. *effektiv*) entscheiden kann. Aus dieser Beschreibung hätte man jeweils ein Modula-, Pascal-, etc. -Programm zur Entscheidung des Problems gewinnen können.

Beim Nachweis der Nichtentscheidbarkeit bzw. Nichtberechenbarkeit ist eine solche intuitive Vorgehensweise nicht mehr möglich, da man hier formal nachweisen muss, dass es *kein Berechnungsverfahren geben kann*. Damit ein solcher Beweis durchführbar ist, benötigt man eine formale Definition dessen, was man unter einem Berechnungsverfahren versteht. Man will also ein *Berechnungsmodell*, das

- 1) **einfach** ist, damit formale Beweise erleichtert werden (z.B. nicht Programmiersprache ADA),
- 2) **berechnungsuniversell** ist, d.h. alle intuitiv berechenbaren Funktionen damit berechnet werden können (z.B. nicht nur endliche Automaten).

Wir werden in diesem Teil drei derartige Modelle betrachten:

- Turingmaschinen
- μ -rekursive Funktionen
- WHILE-Programme

Es gibt noch eine Vielzahl anderer Modelle:

- Minskymaschinen
- GOTO-Programme
- URMs (unbeschränkte Registermaschinen)
- Pascal-Programme
- ...

Es hat sich herausgestellt, dass all diese Modelle *äquivalent* sind, d.h. die gleiche Klasse von Funktionen berechnen. Außerdem ist es bisher nicht gelungen, ein formales Berechnungsmodell zu finden, so dass

- die dadurch berechneten Funktionen noch intuitiv berechenbar erscheinen,
- dadurch Funktionen berechnet werden können, die nicht in den oben genannten Modellen ebenfalls berechenbar sind.

Aus diesen beiden Gründen geht man davon aus, dass die genannten Modelle genau den intuitiven Berechenbarkeitsbegriff formalisieren. Diese Überzeugung nennt man die:

Churchsche These:

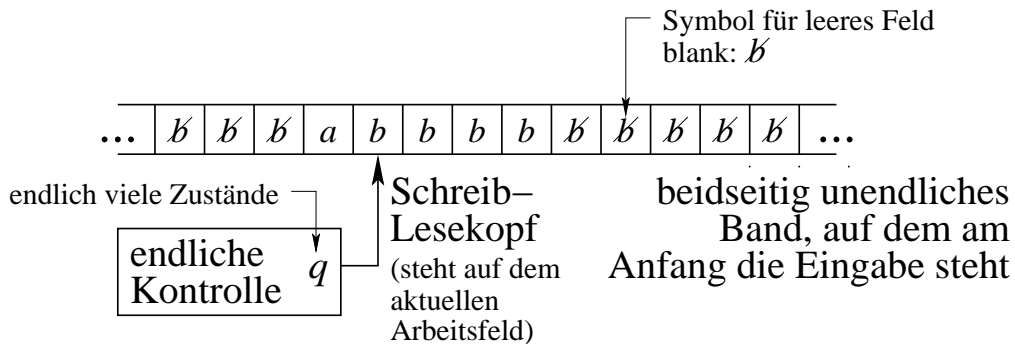
Die (intuitiv) berechenbaren Funktionen sind genau die mit Turingmaschinen (und damit mit GOTO-, WHILE-, Pascal-Programmen, URMs, ...) berechenbaren Funktionen.

Man spricht hier von einer *These* und nicht von einem *Satz*, da es nicht möglich ist, diese Aussage formal zu beweisen. Dies liegt daran, dass der intuitive Berechenbarkeitsbegriff ja nicht formal definierbar ist. Es gibt aber sehr viele Indizien, die für die Richtigkeit der These sprechen (Vielzahl äquivalenter Berechnungsmodelle).

Im folgenden betrachten wir:

- Turingmaschinen
- Zusammenhang zwischen Turingmaschinen und Grammatiken
- Berechenbarkeit, Entscheidbarkeit, Aufzählbarkeit und Zusammenhänge
- Primitiv rekursive Funktionen und LOOP-Programme
- μ -rekursive Funktionen und WHILE-Programme
- Universelle Turingmaschinen und unentscheidbare Probleme
- Weitere unentscheidbare Probleme

11. Turingmaschinen



Zu jedem Zeitpunkt sind nur *endlich viele* Symbole auf dem Band verschieden von \flat .

Definition 11.1 (Turingmaschine)

Eine *Turingmaschine* über dem Eingabealphabet Σ hat die Form $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, \Delta, F)$, wobei

- Q endliche *Zustandsmenge* ist,
- Σ das *Eingabealphabet* ist,
- Γ das *Arbeitsalphabet* ist mit $\Sigma \subseteq \Gamma$, $\flat \in \Gamma \setminus \Sigma$,
- $q_0 \in Q$ der *Anfangszustand* ist,
- $F \subseteq Q$ die *Endzustandsmenge* ist und
- $\Delta \subseteq Q \times \Gamma \times \Gamma \times \{r, l, n\} \times Q$ die *Übergangsrelation* ist.

Dabei bedeutet $(q, a, a', \overset{r}{l}, \underset{n}{q'})$:

- Im Zustand q
- mit a auf dem gerade gelesenen Feld (Arbeitsfeld)

kann die Turingmaschine \mathcal{A}

- das Symbol a durch a' ersetzen,
- in den Zustand q' gehen und
- den Schreib-Lesekopf entweder um ein Feld nach rechts (r), links (l) oder nicht (n) bewegen.

Die Maschine \mathcal{A} heißt *deterministisch*, falls es für jedes Tupel $(q, a) \in Q \times \Gamma$ höchstens ein Tupel der Form $(q, a, \dots, \dots) \in \Delta$ gibt.

NTM steht im folgenden für (möglicherweise nichtdeterministische) Turingmaschinen und *DTM* für deterministische.

Einen Berechnungszustand (*Konfiguration*) einer Turingmaschine kann man beschreiben durch ein Wort $\alpha q \beta$ mit $\alpha, \beta \in \Gamma^+$, $q \in Q$:

- q ist der momentane Zustand
- α ist die Beschriftung des Bandes links vom Arbeitsfeld
- β ist die Beschriftung des Bandes beginnend beim Arbeitsfeld nach rechts

Dabei werden (um endliche Wörter α , β zu erhalten) unendlich viele Blanks weggelassen, d.h. α und β umfassen mindestens den Bandabschnitt, auf dem Symbole $\neq \flat$ stehen.

Beispiel:

Der Zustand der Maschine zu Beginn des Abschnitts wird durch die Konfiguration $aqbbbb$, aber auch durch $\flat baqbbbb \flat$ beschrieben.

Die Übergangsrelation Δ ermöglicht die folgenden *Konfigurationsübergänge*:

Es seien $\alpha, \beta \in \Gamma^+$, $\beta' \in \Gamma^*$, $a, b, a' \in \Gamma$, $q, q' \in Q$.

- $\left. \begin{array}{l} \alpha q a \beta \vdash_{\mathcal{A}} \alpha a' q' \beta \\ \alpha q a \vdash_{\mathcal{A}} \alpha a' q' \flat \end{array} \right\} \text{ falls } (q, a, a', r, q') \in \Delta$
- $\left. \begin{array}{l} \alpha b q a \beta \vdash_{\mathcal{A}} \alpha q' b a' \beta \\ b q a \beta \vdash_{\mathcal{A}} \flat q' b a' \beta \end{array} \right\} \text{ falls } (q, a, a', l, q') \in \Delta$
- $\alpha q a \beta \vdash_{\mathcal{A}} \alpha q' a' \beta \quad \text{falls } (q, a, a', n, q') \in \Delta$

Weitere Bezeichnungen:

- Gilt $k \vdash_{\mathcal{A}} k'$, so heißt k' *Folgekonfiguration* von k .
- Die Konfiguration $\alpha q \beta$ heißt *akzeptierend*, falls $q \in F$.
- Die Konfiguration $\alpha q \beta$ heißt *Stoppkonfiguration*, falls sie keine Folgekonfiguration hat.
- Die von \mathcal{A} akzeptierte Sprache ist
 $L(\mathcal{A}) = \{w \in \Sigma^* \mid \flat q_0 w \flat \vdash_{\mathcal{A}}^* k, \text{ wobei } k \text{ akzeptierende Stoppkonfiguration ist}\}$.

Definition 11.2 (Turing-akzeptierbar, Turing-berechenbar)

- 1) Die Sprache $L \subseteq \Sigma^*$ heißt *Turing-akzeptierbar*, falls es eine NTM \mathcal{A} gibt mit $L = L(\mathcal{A})$.
- 2) Die (partielle) Funktion $f : (\Sigma^*)^n \rightarrow \Sigma^*$ heißt *Turing-berechenbar*, falls es eine DTM \mathcal{A} gibt mit
 - $\forall x_1, \dots, x_n \in \Sigma^* : \flat q_0 x_1 \flat \dots \flat x_n \flat \vdash_{\mathcal{A}}^* k$ mit k Stoppkonfiguration gdw. $(x_1, \dots, x_n) \in \text{dom}(f)$ (Definitionsbereich von f).
 - Im Fall $(x_1, \dots, x_n) \in \text{dom}(f)$ muss die Beschriftung des Bandes in der Stoppkonfiguration k rechts vom Schreib-Lesekopf bis zum ersten Symbol $\notin \Sigma$ der Wert $y = f(x_1, \dots, x_n)$ der Funktion sein, d.h. k muss die Form $uqyv$ haben mit

– $v \in (\Gamma \setminus \Sigma) \cdot \Gamma^* \cup \{\varepsilon\}$

– $y = f(x_1, \dots, x_n)$

Beachte:

- 1) *Undefiniertheit* des Funktionswertes von f entspricht der Tatsache, dass die Maschine bei dieser Eingabe nicht terminiert.
- 2) Bei berechenbaren Funktionen betrachten wir nur *deterministische* Maschinen, da sonst der Funktionswert nicht eindeutig sein müsste.
- 3) Bei $|\Sigma| = 1$ kann man Funktionen von $(\Sigma^*)^n \rightarrow \Sigma^*$ als Funktionen von $\mathbb{N}^k \rightarrow \mathbb{N}$ auffassen (a^k entspricht k).

Beispiel:

Die Funktion

$$f : \mathbb{N} \rightarrow \mathbb{N}, \quad n \mapsto 2n$$

ist Turing-berechenbar. Wie kann eine Turingmaschine die Anzahl der a 's auf dem Band verdoppeln?

Idee:

- Ersetze das erste a durch b ,
- laufe nach rechts bis zum ersten blank, ersetze dieses durch c ,
- laufe zurück bis zum zweiten a (unmittelbar rechts vom b), ersetze dieses durch b ,
- laufe nach rechts bis zum ersten blank etc.
- Sind alle a 's aufgebraucht, so ersetze noch die b 's und c 's wieder durch a 's.

Dies wird durch die folgende *Übergangstafel* realisiert:

$(q_0, \ \not\!b, \ \not\!b, \ n, \ \text{stop}),$	$2 \cdot 0 = 0$
$(q_0, \ a, \ b, \ r, \ q_1),$	ersetze a durch b (\star)
$(q_1, \ a, \ a, \ r, \ q_1),$	laufe nach rechts über a 's
$(q_1, \ c, \ c, \ r, \ q_1),$	und bereits geschriebene c 's
$(q_1, \ \not\!b, \ c, \ n, \ q_2),$	schreibe weiteres c
$(q_2, \ c, \ c, \ l, \ q_2),$	laufe zurück über c 's und
$(q_2, \ a, \ a, \ l, \ q_2),$	a 's
$(q_2, \ b, \ b, \ r, \ q_0),$	bei erstem b eins nach rechts und weiter wie (\star) oder
$(q_0, \ c, \ c, \ r, \ q_3),$	alle a 's bereits ersetzt
$(q_3, \ c, \ c, \ r, \ q_3),$	laufe nach rechts bis Ende der c 's
$(q_3, \ \not\!b, \ \not\!b, \ l, \ q_4),$	letztes c erreicht
$(q_4, \ c, \ a, \ l, \ q_4),$	ersetze c 's und b 's
$(q_4, \ b, \ a, \ l, \ q_4),$	durch a 's
$(q_4, \ \not\!b, \ \not\!b, \ r, \ \text{stop})$	bleibe am Anfang der erzeugten $2n$ a 's stehen

Man sieht hier, dass das Programmieren von Turingmaschinen sehr umständlich ist. Wie bereits erwähnt, betrachtet man solche einfachen (und unpraktischen) Modelle, um das Führen von Beweisen zu erleichtern. Wir werden im folgenden häufig nur die Arbeitsweise einer Turingmaschine beschreiben, ohne die Übergangstafel voll anzugeben.

Beispiel:

Die Sprache $L = \{a^n b^n c^n \mid n \geq 0\}$ ist Turing-akzeptierbar. Die Turingmaschine, welche L akzeptiert, geht wie folgt vor:

- Sie ersetzt das erste a durch a' , das erste b durch b' und das erste c durch c' ;
- läuft zurück zum zweiten a , ersetzt es durch a' , das zweite b durch b' und das zweite c durch c' etc.
- Dies wird solange gemacht, bis nach erzeugtem c' ein $\#$ steht.
- Danach wird von rechts nach links geprüft, ob – in dieser Reihenfolge – nur noch ein c' -Block, dann ein b' -Block und dann ein a' -Block (abgeschlossen durch $\#$) vorhanden ist.
- Die Maschine blockiert, falls dies nicht so ist. Die kann auch bereits vorher geschehen, wenn erwartetes a , b oder c nicht gefunden wird.

Eine entsprechende Turingmaschine \mathcal{A} kann z.B. wie folgt definiert werden:

$$\mathcal{A} = (\{q_0, q_{akz}, finde_b, finde_c, zu_Ende_?, zu_Ende_!, zurück\}, \{a, b, c\}, \{a, a', b, b', c, c', \#\}, q_0, \Delta, \{q_{akz}\}) \text{ mit } \Delta =$$

$(q_0,$	$\#$	$\#$	$N,$	$q_{akz}),$
$(q_0,$	$a,$	$a',$	$R,$	$finde_b),$
$(finde_b,$	$a,$	$a,$	$R,$	$finde_b),$
$(finde_b,$	$b',$	$b',$	$R,$	$finde_b),$
$(finde_b,$	$b,$	$b',$	$R,$	$finde_c),$
$(finde_c,$	$b,$	$b,$	$R,$	$finde_c),$
$(finde_c,$	$c',$	$c',$	$R,$	$finde_c),$
$(finde_c,$	$c,$	$c',$	$R,$	$zu_Ende_?),$
$(zu_Ende_?,$	$c,$	$c,$	$L,$	$zurück),$
$(zurück,$	$c',$	$c',$	$L,$	$zurück),$
$(zurück,$	$b,$	$b,$	$L,$	$zurück),$
$(zurück,$	$b',$	$b',$	$L,$	$zurück),$
$(zurück,$	$a,$	$a,$	$L,$	$zurück),$
$(zurück,$	$a',$	$a',$	$R,$	$q_0),$
$(zu_Ende_?,$	$\#$	$\#$	L	$zu_Ende_!),$
$(zu_Ende_!,$	$c',$	$c',$	$L,$	$zu_Ende_!),$
$(zu_Ende_!,$	$b',$	$b',$	$L,$	$zu_Ende_!),$
$(zu_Ende_!,$	$a',$	$a',$	$L,$	$zu_Ende_!),$
$(zu_Ende_!,$	$\#$	$\#$	$N,$	$q_{akz})\}$

Varianten von Turingmaschinen:

In der Literatur werden verschiedene Versionen der Definition der Turingmaschine angegeben, die aber alle äquivalent zueinander sind, d.h. dieselben Sprachen akzeptieren und dieselben Funktionen berechnen. Hier zwei Beispiele:

- Turingmaschinen mit nach links begrenztem und nur nach rechts unendlichem Arbeitsband
- Turingmaschinen mit mehreren Bändern und Schreib-Leseköpfen

Wir betrachten das zweite Beispiel genauer und zeigen Äquivalenz zur in Definition 11.1 eingeführten 1-Band-TM.

Definition 11.3 (*k*-Band-TM)

Eine *k*-Band-NTM hat die Form $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, \Delta, F)$ mit

- $Q, \Sigma, \Gamma, q_0, F$ wie in Definition 11.1 und
- $\Delta \subseteq Q \times \Gamma^k \times \Gamma^k \times \{r, l, n\}^k \times Q$.

Dabei bedeutet $(q, (a_1, \dots, a_k), (b_1, \dots, b_k), (d_1, \dots, d_k), q') \in \Delta$:

- Vom Zustand q aus
- mit a_1, \dots, a_k auf den Arbeitsfeldern der k Bänder

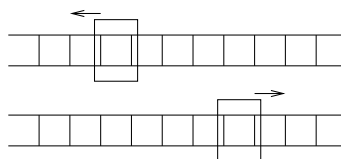
kann \mathcal{A}

- das Symbol a_i auf dem i -ten Band durch b_i ersetzen,
- in den Zustand q' gehen und
- die Schreib-Leseköpfe der Bänder entsprechend d_i bewegen.

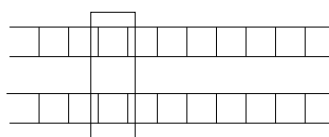
Das erste Band wird (o.B.d.A.) als Ein- und Ausgabeband verwendet.

Beachte:

Wichtig ist hier, dass sich die Köpfe der verschiedenen Bänder auch verschieden bewegen können:



Wären die Köpfe gekoppelt, so hätte man im Prinzip nicht mehrere Bänder, sondern ein Band mit mehreren Spuren:



k Spuren erhält man einfach, indem man eine normale NTM (nach Definition 11.1) verwendet, die als Bandalphabet Γ^k statt Γ hat.

Offenbar kann man jede 1-Band-NTM (nach Definition 11.1) durch eine k -Band-NTM ($k > 1$) simulieren, indem man nur das erste Band wirklich verwendet. Der nächste Satz zeigt, dass auch die Umkehrung gilt:

Satz 11.4

Wird die Sprache L durch eine k -Band-NTM akzeptiert, so auch durch eine 1-Band-NTM.

Beweis. Es sei \mathcal{A} eine k -Band-NTM. Gesucht ist

- eine 1-Band-NTM \mathcal{A}' und
- eine Kodierung $k \mapsto k'$ der Konfigurationen k von \mathcal{A} durch Konfigurationen k' von \mathcal{A}' ,

so dass gilt:

$$k_1 \vdash_{\mathcal{A}} k_2 \text{ gdw. } k'_1 \vdash_{\mathcal{A}'}^* k'_2$$

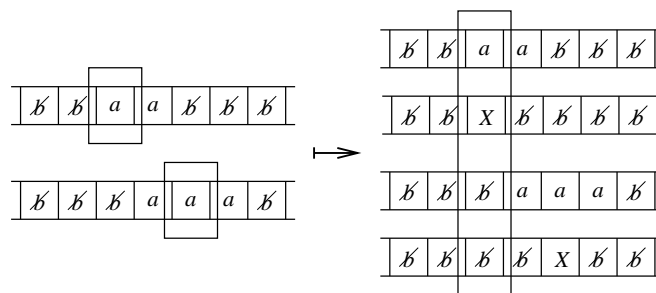
(Zur Simulation eines Schrittes von \mathcal{A} sind mehr Schritte nötig).

Arbeitsalphabet von \mathcal{A}' : $\Gamma^{2k} \cup \Sigma \cup \{\beta\}$;

- $\Sigma \cup \{\beta\}$ wird für Eingabe benötigt;
- Γ^{2k} sorgt dafür, dass sich \mathcal{A}' wie eine 1-Band-NTM mit $2k$ Spuren verhält.

Kodierung: Jeweils 2 Spuren kodieren ein Band von \mathcal{A} .

- Die erste Spur enthält die Bandbeschriftung.
- Die zweite Spur enthält eine Markierung X (und sonst Blanks), die zeigt, wo das Arbeitsfeld des Bandes ist, z.B.



Hierbei wird angenommen, dass das Arbeitsfeld von \mathcal{A}' stets bei dem am weitesten links stehenden X liegt.

Initialisierung: Zunächst wird die Anfangskonfiguration

$$\#q_0a_1 \dots a_m\#$$

von \mathcal{A}' in die *Kodierung* der entsprechenden Anfangskonfiguration von \mathcal{A} umgewandelt (durch entsprechende Übergänge):

$\#$	$\#$	a_1	a_2	\dots	a_m	$\#$	Spur 1 und 2 kodieren
$\#$	$\#$	X	$\#$	\dots	$\#$	$\#$	Band 1
$\#$	$\#$	$\#$	$\#$	\dots	$\#$	$\#$	Spur 3 und 4 kodieren
$\#$	$\#$	X	$\#$	\dots	$\#$	$\#$	Band 2

⋮

Simulation der Übergänge: Betrachte $(q, (a_1, \dots, a_k), (b_1, \dots, b_k), (d_1, \dots, d_k), q') \in \Delta$.

- Von links nach rechts suche die mit X markierten Felder. Dabei merke man sich (in dem Zustand der TM) die Symbole, welche jeweils über dem X stehen. Außerdem zählt man (durch Zustand der TM) mit, wie viele X man schon gelesen hat, um festzustellen, wann das k -te erreicht ist. Man bestimmt so, ob das aktuelle Tupel tatsächlich (a_1, \dots, a_k) ist.
- Von rechts nach links gehend überdrucke man die a_i bei den X -Marken jeweils durch das entsprechende b_i und verschiebt die X -Marken gemäß d_i .
- Bleibe bei der am weitesten links stehenden X -Markierung, lasse den Kopf dort stehen und gehe in Zustand q' .

□

Bemerkung 11.5.

- 1) War \mathcal{A} deterministisch, so liefert obige Konstruktion auch eine deterministische 1-Band-Turingmaschine.
- 2) Diese Konstruktion kann auch verwendet werden, wenn man sich für die berechnete Funktion interessiert. Dazu muss man am Schluss (wenn \mathcal{A} in Stoppkonfiguration ist) in der Maschine \mathcal{A}' noch die Ausgabe geeignet aufbereiten.

Bei der Definition von Turing-berechenbar haben wir uns von vornherein auf *deterministische* Turingmaschinen beschränkt. Der folgende Satz zeigt, dass man dies auch bei Turing-akzeptierbaren Sprachen machen kann, ohne an Ausdruckstärke zu verlieren.

Satz 11.6

Zu jeder NTM gibt es eine DTM, die dieselbe Sprache akzeptiert.

Beweis. Wegen Satz 11.4 und Bemerkung 11.5 genügt es, eine deterministische 3-Band-Turingmaschine zu konstruieren. Es sei $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, \Delta, F)$ eine NTM.

Die Maschine \mathcal{A}' soll für wachsendes n auf dem dritten Band jeweils alle Konfigurationsfolgen

$$k_0 \vdash_{\mathcal{A}} k_1 \vdash_{\mathcal{A}} k_2 \vdash_{\mathcal{A}} \dots \vdash_{\mathcal{A}} k_n$$

beginnend mit der Startkonfiguration $k_0 = \#q_0w\#$ erzeugen.

Die Kontrolle, dass tatsächlich alle solchen Folgen erzeugt werden, wird auf dem zweiten Band vorgenommen. Das erste Band speichert das Eingabewort w (damit man stets weiß, was k_0 sein muss).

Genauer: Es sei

$$r = \text{maximale Anzahl von Transitionen in } \Delta \text{ pro festem Paar } (q, a) \in Q \times \Gamma$$

(entspricht dem maximalen Verzweigungsgrad der nichtdeterministischen Berechnung).

Eine Indexfolge i_1, \dots, i_n mit $i_j \in \{1, \dots, r\}$ bestimmt dann von k_0 aus für n Schritte die Auswahl der jeweiligen Transition, und somit von k_0 aus eine feste Konfigurationsfolge

$$k_0 \vdash_{\mathcal{A}} k_1 \vdash_{\mathcal{A}} \dots \vdash_{\mathcal{A}} k_n$$

Zählt man daher alle endlichen Wörter über $\{1, \dots, r\}$ auf und zu jedem Wort $i_1 \dots i_n$ die zugehörige Konfigurationsfolge, so erhält man eine Aufzählung aller endlichen Konfigurationsfolgen.

\mathcal{A}' realisiert dies auf den drei Bändern wie folgt:

- Auf Band 1 bleibt die Eingabe gespeichert.
- Auf dem zweiten Band werden sukzessive alle Wörter $i_1 \dots i_n \in \{1, \dots, r\}^*$ erzeugt.
- Für jedes dieser Wörter wird auf dem dritten Band die zugehörige Konfigurationsfolge realisiert. Erreicht man hierbei eine akzeptierende Konfiguration von \mathcal{A} , so geht auch \mathcal{A}' in eine akzeptierende Konfiguration.

□

12. Zusammenhang zwischen Turingmaschinen und Grammatiken

Wir zeigen zunächst den Zusammenhang zwischen den *Typ-0-Sprachen* und *Turing-akzeptierbaren Sprachen*. Dieser beruht darauf, dass es eine Entsprechung von Ableitungen einer Typ-0-Grammatik einerseits und Konfigurationsfolgen von Turingmaschinen andererseits gibt. Beim Übergang von der Turingmaschine zur Grammatik dreht sich allerdings die Richtung um:

- eine akzeptierende Konfigurationsfolge beginnt mit dem zu akzeptierenden Wort
- eine Ableitung endet mit dem erzeugten Wort

Satz 12.1

Eine Sprache L gehört zu \mathcal{L}_0 gdw. sie Turing-akzeptierbar ist.

Beweis.

„ \Rightarrow “: Es sei $L = L(G)$ für eine Typ-0-Grammatik $G = (N, \Sigma, P, S)$ mit $|P| = k$ Regeln. Wir geben eine 2-Band-NTM an, die $L(G)$ akzeptiert (die nach Satz 11.4 äquivalent zu einer 1-Band-NTM ist).

1. Band: speichert Eingabe w

2. Band: ausgehend von S werden gemäß den Regeln von G ableitbare Wörter gebildet

Es wird verglichen, ob auf Band 2 irgendwann w (d.h. der Inhalt von Band 1) entsteht. Wenn ja, so geht man in akzeptierende Stoppkonfiguration, sonst sucht die Turingmaschine weiter.

Die Maschine geht dabei wie folgt vor:

- 1) Schreibe S auf Band 2, gehe nach links auf das $\$$ vor S .
- 2) Gehe auf Band 2 nach rechts und wähle (nichtdeterministisch) eine Stelle aus, an der die linke Seite der anzuwendenden Produktion beginnen soll.
- 3) Wechsle nun (nichtdeterministisch) in einen der Zustände $\text{Regel}_1, \dots, \text{Regel}_k$ (entspricht der Entscheidung, dass man diese Regel anwenden will).

- 4) Es sei Regel_i der gewählte Zustand und $\alpha \rightarrow \beta$ die entsprechende Produktion.

Überprüfe, ob α tatsächlich die Beschriftung des Bandstücks der Länge $|\alpha|$ ab der gewählten Bandstelle ist.

- 5) Falls der Test erfolgreich war, so ersetze α durch β .

Vorher müssen bei $|\alpha| < |\beta|$ die Symbole $\neq \$$ rechts von α um $|\beta| - |\alpha|$ Positionen nach rechts

bzw. bei $|\alpha| > |\beta|$ um $|\alpha| - |\beta|$ Positionen nach links geschoben werden.

6) Gehe nach links bis zum ersten $\#$ und vergleiche, ob die Beschriftung auf dem Band 1 mit der auf Band 2 übereinstimmt.

7) Wenn ja, so gehe in akzeptierenden Stoppzustand. Sonst fahre fort bei 2).

„ \Leftarrow “: Es sei $L = L(\mathcal{A})$ für eine NTM $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, \Delta, F)$.

Wir konstruieren eine Grammatik G , die jedes Wort $w \in L(\mathcal{A})$ wie folgt erzeugt:

1. Phase: Erst wird w mit „genügend vielen“ $\#$ -Symbolen links und rechts davon erzeugt (dies passiert für jedes w , auch für $w \notin L(\mathcal{A})$).

„Genügend viele“ bedeutet dabei: so viele, wie \mathcal{A} beim Akzeptieren von w vom Arbeitsband benötigt.

2. Phase: Auf dem so erzeugten Arbeitsband simuliert G die Berechnung von \mathcal{A} bei Eingabe w .

3. Phase: War die Berechnung akzeptierend, so erzeuge nochmals das w .

Damit man in der zweiten Phase das in der dritten Phase benötigte w nicht vergisst, verwendet man als Nichtterminalsymbole Tupel aus

$$(\Sigma \cup \{\#\}) \times \Gamma$$

wobei man

- sich in der ersten Komponente w merkt und
- in der zweiten Komponente die Berechnung simuliert.

Formale Definition:

$$N = \{S, A, B, E\} \cup Q \cup (\Sigma \cup \{\#\}) \times \Gamma,$$

wobei

- S, A, B zum Aufbau des Rechenbandes am Anfang,
- E zum Löschen am Schluss,
- Q zur Darstellung des aktuellen Zustandes,
- $\Sigma \cup \{\#\}$ zum Speichern von w und
- Γ zur \mathcal{A} -Berechnung

dienen.

Regeln:

1. Phase: Erzeuge w und ein genügend großes Arbeitsband.

$$S \longrightarrow Bq_0A$$

$$A \longrightarrow [a, a]A \text{ für alle } a \in \Sigma$$

$$A \longrightarrow B$$

$$B \longrightarrow [\#, \#]B$$

$$B \longrightarrow \varepsilon$$

Man erhält also somit für alle $a_1 \dots a_n \in \Sigma^*$, $k, l, \geq 0$:

$$S \vdash_G^* [\mathcal{B}, \mathcal{B}]^k q_0[a_1, a_1] \dots [a_n, a_n][\mathcal{B}, \mathcal{B}]^l$$

2. Phase: simuliert TM-Berechnung in der „zweiten Spur“:

- $p[a, b] \longrightarrow [a, b']q$
falls $(p, b, b', r, q) \in \Delta$, $a \in \Sigma \cup \{\mathcal{B}\}$
- $[a, c]p[a', b] \longrightarrow q[a, c][a', b']$
falls $(p, b, b', l, q) \in \Delta$, $a, a' \in \Sigma \cup \{\mathcal{B}\}$, $c \in \Gamma$
- $p[a, b] \longrightarrow q[a, b']$
falls $(p, b, b', n, q) \in \Delta$, $a \in \Sigma \cup \{\mathcal{B}\}$

Beachte:

Da wir in der ersten Phase genügend viele Blanks links und rechts von $a_1 \dots a_n$ erzeugen können, muss in der zweiten Phase das „Nachschieben“ von Blanks am Rand nicht mehr behandelt werden.

3. Phase: Aufräumen und erzeugen von $a_1 \dots a_n$, wenn die TM akzeptiert hat

- $q[a, b] \longrightarrow EaE$ für $a \in \Sigma$, $b \in \Gamma$
- $q[\mathcal{B}, b] \longrightarrow E$ für $b \in \Gamma$
falls $q \in F$ und es keine Transition der Form $(q, b, \dots, \dots) \in \Delta$ gibt (d.h. akzeptierende Stoppkonfiguration erreicht)
- $E[a, b] \longrightarrow aE$ für $a \in \Sigma$, $b \in \Gamma$
(Aufräumen nach rechts)
- $[a, b]E \longrightarrow Ea$ für $a \in \Sigma$, $b \in \Gamma$
(Aufräumen nach links)
- $E[\mathcal{B}, b] \longrightarrow E$ für $b \in \Gamma$
(Entfernen des zusätzlich benötigten Arbeitsbandes nach rechts)
- $[\mathcal{B}, b]E \longrightarrow E$ für $b \in \Gamma$
(Entfernen des zusätzlich benötigten Arbeitsbandes nach links)
- $E \longrightarrow \varepsilon$

Man sieht nun leicht, dass für alle $w \in \Sigma^*$ gilt:

$$w \in L(G) \text{ gdw. } \mathcal{A} \text{ akzeptiert } w. \quad \square$$

Für Typ-0-Sprachen gelten die folgenden Abschlusseigenschaften:

Satz 12.2

- 1) \mathcal{L}_0 ist abgeschlossen unter $\cup, \cdot, *$ und \cap .

2) \mathcal{L}_0 ist nicht abgeschlossen unter Komplement.

Beweis.

1) Für die regulären Operationen $\cup, \cdot, *$ zeigt man dies im Prinzip wie für \mathcal{L}_2 durch Konstruktion einer entsprechenden Grammatik.

Damit sich die Produktionen der verschiedenen Grammatiken nicht gegenseitig beeinflussen, genügt es allerdings nicht mehr, nur die Nichtterminalsymbole der Grammatiken disjunkt zu machen.

Zusätzlich muss man die Grammatiken in die folgende Form bringen:

Die Produktionen sind von der Form

$$\begin{array}{ll} u \longrightarrow v & \text{mit } u \in N_i^+ \text{ und } v \in N_i^* \\ X_a \longrightarrow a & \text{mit } X_a \in N_i \text{ und } a \in \Sigma \end{array}$$

Für den Durchschnitt verwendet man Turingmaschinen:

Die NTM für $L_1 \cap L_2$ verwendet zwei Bänder und simuliert zunächst auf dem ersten die Berechnung der NTM für L_1 und dann auf dem anderen die der NTM für L_2 .

Wenn beide zu akzeptierenden Stoppkonfigurationen der jeweiligen Turingmaschinen führen, so geht die NTM für $L_1 \cap L_2$ in eine akzeptierende Stoppkonfiguration.

Beachte:

Es kann sein, dass die NTM für L_1 auf einer Eingabe w nicht terminiert, die NTM für $L_1 \cap L_2$ also gar nicht dazu kommt, die Berechnung der NTM für L_2 zu simulieren. Aber dann ist ja w auch nicht in $L_1 \cap L_2$.

2) Wir werden später sehen, dass Turing-akzeptierbare Sprachen nicht unter Komplement abgeschlossen sind (Satz 16.10). □

Wir werden später außerdem zeigen, dass für Turing-akzeptierbare Sprachen (und damit für \mathcal{L}_0) alle bisher betrachteten Entscheidungsprobleme unentscheidbar sind (Sätze 16.6, 16.8, 16.9):

Satz 12.3

Für \mathcal{L}_0 sind das Leerheitsproblem, das Wortproblem und das Äquivalenzproblem unentscheidbar.

Von den Sprachklassen aus der Chomsky-Hierarchie sind nun alle bis auf \mathcal{L}_1 (kontextsensitiv) durch geeignete Automaten/Maschinenmodelle charakterisiert. Man kann mit der Beweisidee von Satz 12.1 auch eine Charakterisierung kontextsensitiver Sprachen durch spezielle Turingmaschinen erhalten.

Diese Maschinen dürfen nur auf dem Bandabschnitt arbeiten, auf dem anfangs die Eingabe stand. Um ein Überschreiten der dadurch gegebenen Bandgrenzen zu verhindern, verwendet man Randmarker $\phi, \$$.

Definition 12.4 (linear beschränkter Automat)

Ein *linear beschränkter Automat (LBA)* ist eine NTM $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, \Delta, F)$, so dass

- $\$, \not\in \Gamma \setminus \Sigma$
- Übergänge (q, \not, \dots) sind nur in der Form (q, \not, \not, r, q') erlaubt (linker Rand darf nicht überschritten werden).
- Übergänge $(q, \$, \dots)$ sind nur in der Form $(q, \$, \$, l, q')$ erlaubt.
- \not und $\$$ dürfen nicht geschrieben werden.

Ein gegebener LBA \mathcal{A} akzeptiert die Sprache

$$L(\mathcal{A}) := \{w \in \Sigma^* \mid \not q_0 w \$ \vdash^* k, \text{ wobei } k \text{ akzeptierende Stoppkonfiguration ist}\}.$$

Korollar 12.5

Eine Sprache L gehört zu \mathcal{L}_1 gdw. sie von einem LBA akzeptiert wird.

Beweis.

„ \Rightarrow “: Verwende die Konstruktion aus dem Beweis zu Satz 12.1.

Da alle Produktionen von kontextsensitiven Grammatiken nichtkürzend sind (mit Ausnahme $S \rightarrow \varepsilon$), muss man auf dem zweiten Band nur ableitbare Wörter bis zur Länge $|w|$ erzeugen (längere können nie mehr zu w abgeleitet werden). Daher kommt man mit $|w|$ vielen Feldern aus.

Beachte:

Zwei Bänder liefern nicht ein doppelt so langes Band, sondern ein größeres Alphabet.

„ \Leftarrow “: Durch Modifikation der Beweisidee von Satz 12.1 gelingt es, zu einem LBA eine Grammatik zu konstruieren, die nur nichtkürzende Regeln hat

$$\alpha \longrightarrow \beta \text{ mit } |\alpha| \leq |\beta|$$

Wie im Beweis von Satz 9.3 illustriert, kann man diese in eine kontextsensitive Grammatik transformieren.

Idee:

Da man mit $|w|$ Arbeitsfeldern auskommt, muss man keine $[\not, \not]$ links und rechts von w erzeugen. Dadurch fallen dann auch die folgenden kürzenden Regeln weg:

$$E[\not, \not] \longrightarrow E$$

$$[\not, \not]E \longrightarrow E$$

Es gibt allerdings noch einige technische Probleme:

- Man muss die Randmarker \not und $\$$ einführen und am Schluss löschen.

- Man muss das Hilfssymbol E und den Zustand q löschen.

Lösung:

Führe die Symbole ℓ , $\$$ und den Zustand q nicht als zusätzliche Symbole ein, sondern kodiere sie in die anderen Symbole hinein.

z.B. statt $[a, b]q[a', b'][a'', b'']$ verwende $[a, b][q, a', b'][a'', b'']$.

□

Satz 12.6

\mathcal{L}_1 ist abgeschlossen unter $\cup, \cdot, *, \cap$ und Komplement.

Beweis. Für $\cup, \cdot, *$ und \cap geht dies wie bei \mathcal{L}_0 .

Komplement: schwierig, war lange offen und wurde dann in den 1980ern unabhängig von zwei Forschern gezeigt (Immerman und Szelepcsényi). \square

Für LBAs ist bisher nicht bekannt, ob deterministische LBAs genauso stark wie nicht-deterministische LBAs sind.

Satz 12.7

Für \mathcal{L}_1 ist das Wortproblem entscheidbar.

Beweis. Da kontextsensitive Produktionen (bis auf Spezialfall $S \rightarrow \varepsilon$) nichtkürzend sind, muss man zur Entscheidung „ $w \in L(G)$?“ nur alle aus S ableitbaren Wörter aus $(N \cup \Sigma)^*$ der Länge $\leq |w|$ erzeugen.

Dies sind endlich viele. \square

13. Berechenbarkeit, Entscheidbarkeit, Aufzählbarkeit

Wir werden hier – gemäß der These von Church – die Begriffe *Turing-berechenbar* und (*intuitiv*) *berechenbar* als *synonym* verwenden.

Wegen der anfangs erwähnten Äquivalenz von Berechnungsmodellen treffen alle Definitionen und Resultate auch auf die anderen Modelle zu.

Wir werden meist intuitiv die Existenz eines Berechnungsverfahrens (einer DTM) begründen. Diese intuitiven Argumente können aber leicht (wenn auch im Detail technisch und zeitaufwendig) in TM-Konstruktionen übersetzt werden.

Definition 13.1 (Berechenbarkeit)

Wir betrachten partielle oder totale Funktionen

$$f : (\Sigma^*)^n \rightarrow \Sigma^*.$$

- f heißt *berechenbar* (*partiell rekursiv*), falls sie Turing-berechenbar ist.
- Ist f total (d.h. $\text{dom}(f) = (\Sigma^*)^n$), so heißt f *rekursiv*.

Bei partiellen Funktionen entspricht undefiniertheit des Funktionswertes der Tatsache, dass das Berechnungsverfahren (die DTM) bei dieser Eingabe nicht terminiert.

Definition 13.2 (entscheidbar, partiell entscheidbar, rekursiv aufzählbar)

Wir betrachten n -stellige Relationen $R \subseteq (\Sigma^*)^n$.

- 1) R heißt *entscheidbar* (*rekursiv*), falls ihre *charakteristische Funktion*

$$\chi_R : (\Sigma^*)^n \rightarrow \Sigma^* \text{ mit } (x_1, \dots, x_n) \mapsto \begin{cases} a & \text{falls } (x_1, \dots, x_n) \in R \\ \varepsilon & \text{sonst} \end{cases}$$

berechenbar ist. Dabei ist a ein Element von Σ (beliebig aber fest gewählt).

- 2) R heißt *partiell entscheidbar* (*partiell rekursiv*), falls R Definitionsbereich einer berechenbaren Funktion f ist.

D.h. es gibt eine DTM, die bei Eingabe $(x_1, \dots, x_n) \in (\Sigma^*)^n$

- terminiert, falls $(x_1, \dots, x_n) \in R$ ist und
- nicht terminiert sonst.

- 3) R heißt *rekursiv aufzählbar*, falls R von einer Aufzähl-Turingmaschine aufgezählt wird. Diese ist wie folgt definiert:

Eine *Aufzähl-Turingmaschine* \mathcal{A} ist eine DTM, die einen speziellen Ausgabezustand q_{Ausgabe} hat.

Eine *Ausgabekonfiguration* für eine n -stellige Relation ist von der Form

$$uq_{\text{Ausgabe}}x_1\#x_2\#\dots\#x_n\#v \text{ mit } u, v \in \Gamma^* \text{ und } x_1, \dots, x_n \in \Sigma^*.$$

Diese Konfiguration hat (x_1, \dots, x_n) als *Ausgabe*.

Die durch \mathcal{A} aufgezählte Relation ist

$$R = \{(x_1, \dots, x_n) \in (\Sigma^*)^n \mid (x_1, \dots, x_n) \text{ ist Ausgabe einer Ausgabekonfiguration, die von } \mathcal{A} \text{ ausgehend von } \#q_0\# \text{ erreicht wird}\}.$$

Bemerkung 13.3.

- 1) Probleme wie z.B. das *Wortproblem* für kontextfreie Sprachen können als *Relationen* aufgefasst werden:

Es sei $G = (N, \Sigma, P, S)$ eine kontextfreie Grammatik und $w \in \Sigma^*$.

Die Grammatik G kann als Wort $\text{code}(G) \in \Gamma^*$ über einem erweiterten Alphabet Γ aufgefasst werden. Das Wortproblem für G entspricht der Relation

$$R = \{(\text{code}(G), w) \mid w \in L(G)\} \subseteq (\Gamma^*)^2.$$

- 2) Für $n = 1$ stimmen die Begriffe „partiell entscheidbar“ und „Turing-akzeptierbar“ überein, denn:

Bei Turing-akzeptierbar war neben dem Anhalten der (o.B.d.A. deterministischen) TM zusätzlich gefordert, dass man einen akzeptierenden Zustand erreicht hat. Man kann aber einfach von allen nichtakzeptierenden Stoppkonfigurationen aus in eine Endlosschleife gehen.

- 3) Im Fall $n = 1$ sind folgende Charakterisierungen äquivalent:

- (a) R ist rekursiv aufzählbar
- (b) $R = \emptyset$ oder R ist Wertebereich einer rekursiven Funktion.
- (c) R ist Wertebereich einer partiell rekursiven Funktion.

Beweis.

„ $a \Rightarrow b$ “ Es sei R rekursiv aufzählbar und $R \neq \emptyset$. Weiterhin sei \mathcal{A} eine Aufzähl-TM für R .

Wähle $w_0 \in R$ beliebig.

Wir definieren die Funktion $f : \Sigma^* \times \Sigma^* \rightarrow \Sigma^*$ wie folgt:

$$f := \begin{cases} u & \text{falls } \mathcal{A} \text{ nach } \leq |V| \text{ Schritte } u \text{ als Ausgabe erzeugt} \\ w_0 & \text{sonst} \end{cases}$$

Es ist leicht zu sehen, dass der Wertebereich $\text{ran}(f)$ von f die Menge R ist.

Außerdem kann man aus \mathcal{A} leicht eine TM erhalten, die f berechnet.

„ $b \Rightarrow c$ “ Für $R \neq \emptyset$ ist dies trivial.

Für $R = \emptyset$ betrachten wir die überall undefinierte Funktion. Diese ist offenbar partiell berechenbar und hat als Wertebereich die leere Menge.

„ $c \Rightarrow a$ “ Es sei $R = \text{ran}(f)$ und \mathcal{A} eine DTM, die f berechnet.

Die Aufzähl-TM \mathcal{B} für R arbeitet wie folgt:

Es sei e_1, e_2, \dots eine effektive Aufzählung aller Eingaben für \mathcal{A} .

Dann lässt \mathcal{B}

- 1) die Berechnung von \mathcal{A} bei Eingabe e_1 *einen Schritt* laufen
- 2) die Berechnung von \mathcal{A} bei Eingabe e_1 *zwei Schritte* und bei Eingabe e_2 ebenfalls *2 Schritte* laufen
- ⋮
- n) die Berechnung von \mathcal{A} mit den Eingaben e_1, e_2, \dots, e_n jeweils *n Schritte* laufen
- ⋮

Wird dabei ein Funktionswert von f berechnet, so wird dieser ausgegeben.

□

Satz 13.4

Es sei $R \subseteq (\Sigma^*)^n$ eine Relation.

- 1) R ist rekursiv aufzählbar gdw. R ist partiell entscheidbar.
- 2) Ist R entscheidbar, so auch partiell entscheidbar.
- 3) Ist R entscheidbar, so auch $\overline{R} = (\Sigma^*)^n \setminus R$.
- 4) R ist entscheidbar gdw. R und \overline{R} partiell entscheidbar sind.

Beweis.

- 1) „ \Rightarrow “: Es sei R rekursiv aufzählbar und \mathcal{A} eine Aufzähl-DTM für R . Die Maschine \mathcal{A}' arbeitet wie folgt:

- Sie speichert die Eingabekonfiguration $\#q_0 \#x_1 \#x_2 \dots \#x_n \#$ ab (z.B. auf zusätzlichem Band).
- Sie beginnt (auf anderem Band) mit der Aufzählung von R .
- Bei jeder Ausgabekonfiguration überprüft sie, ob die entsprechende Ausgabe mit der gespeicherten Eingabe übereinstimmt. Wenn ja, so terminiert \mathcal{A}' . Sonst sucht sie die nächste Ausgabekonfiguration von \mathcal{A} .
- Terminiert \mathcal{A} , ohne dass (x_1, \dots, x_n) ausgegeben wurde, so gehe in Endlosschleife.

\mathcal{A}' terminiert daher genau dann nicht, wenn (x_1, \dots, x_n) nicht in der Aufzählung vorkommt.

- „ \Leftarrow “: Es sei \mathcal{A} ein partielles Entscheidungsverfahren für R , d.h. \mathcal{A} berechnet eine Funktion f mit $\text{dom}(f) = R$.

Es sei $\vec{x}^{(1)}, \vec{x}^{(2)}, \vec{x}^{(3)}, \dots$ eine Auflistung von $(\Sigma^*)^n$.

Die Maschine \mathcal{A}' arbeitet wie folgt:

- 1) Simuliere die Berechnung von \mathcal{A} mit der Eingabe $\vec{x}^{(1)}$ *einen Schritt*.
- 2) Simuliere die Berechnung von \mathcal{A} mit der Eingabe $\vec{x}^{(1)}$ *zwei Schritte* und mit der Eingabe $\vec{x}^{(2)}$ *zwei Schritte*.
- ⋮
- n) Simuliere die Berechnung von \mathcal{A} mit den Eingaben $\vec{x}^{(1)}, \dots, \vec{x}^{(n)}$ jeweils *n Schritte*.
- ⋮

Terminiert \mathcal{A} für eine dieser Eingaben, so gebe diese Eingabe aus und mache weiter.

Beachte:

Man kann nicht \mathcal{A} zunächst ganz mit Eingabe $\vec{x}^{(1)}$ laufen lassen (ohne Schrittbeschränkung), da \mathcal{A} darauf nicht terminieren muss. Das gewählte Vorgehen nennt man *dove-tailing* (Ineinander-Verzahnungen mehrerer Berechnungen).

- 2) Eine DTM \mathcal{A} , die χ_R berechnet, wird wie folgt modifiziert:
 - Ist die Ausgabe ε (d.h. in der Stoppkonfiguration von \mathcal{A} steht der Kopf auf \flat), so gehe in Endlosschleife.
 - Sonst gehe in Stoppkonfiguration.
- 3) Eine DTM \mathcal{A} , die χ_R berechnet, wird wie folgt zu einer DTM für $\chi_{\bar{R}}$ modifiziert:
 - Ist die Ausgabe von \mathcal{A} das leere Wort ε , so erzeuge Ausgabe a .
 - Sonst erzeuge Ausgabe ε .
- 4) „ \Rightarrow “: Ergibt sich aus 2) und 3).

„ \Leftarrow “: Sind R und \bar{R} partiell entscheidbar, so mit 1) auch rekursiv aufzählbar.

Für Eingabe \vec{x} lässt man die Aufzähl-DTMs \mathcal{A} und \mathcal{B} für R und \bar{R} parallel laufen (d.h. jeweils abwechselnd ein Schritt von \mathcal{A} auf einem Band gefolgt von einem Schritt von \mathcal{B} auf dem anderen).

Die Eingabe \vec{x} kommt in einer der beiden Aufzählungen vor:

$$\vec{x} \in (\Sigma^*)^n = R \cup \bar{R}$$

Kommt \vec{x} bei \mathcal{A} vor, so erzeuge Ausgabe a , sonst Ausgabe ε .

□

Für $f : (\Sigma^*)^n \rightarrow \Sigma^*$ definieren wir:

$$\text{Graph}(f) := \{(x_1, \dots, x_n, x_{n+1}) \mid (x_1, \dots, x_n) \in \text{dom}(f) \text{ und } f(x_1, \dots, x_n) = x_{n+1}\}$$

Satz 13.5

f ist berechenbar g.d.w. $\text{Graph}(f)$ ist rekursiv aufzählbar.

Beweis.

„ \Leftarrow “: Sei $\text{Graph}(f)$ rekursiv aufzählbar. Um den Funktionswert $f(x_1, \dots, x_n)$ zu berechnen, startet man das Aufzählungsverfahren und wartet, bis ein Tupel der Form (x_1, \dots, x_n, y) auftaucht.

- Wenn ja, so ist y der Funktionswert.
- Sonst ist $(x_1, \dots, x_n) \notin \text{dom}(f)$, und die Maschine terminiert nicht, d.h. sie zählt unendlich auf oder geht in Endlosschleife, wenn die Aufzählung abgebrochen wird, ohne dass Tupel der Form (x_1, \dots, x_n, \cdot) gefunden wurde).

„ \Rightarrow “: Mit Satz 13.4 genügt es zu zeigen, dass $\text{Graph}(f)$ partiell entscheidbar ist.

Das partielle Entscheidungsverfahren für $\text{Graph}(f)$ erhält man wie folgt:

Berechne bei Eingabe $(x_1, \dots, x_n, x_{n+1})$ den Funktionswert von f an der Stelle (x_1, \dots, x_n) .

1. Fall: $f(x_1, \dots, x_n)$ ist definiert und $= x_{n+1}$.
Dann terminiere.

2. Fall: $f(x_1, \dots, x_n)$ ist definiert und $\neq x_{n+1}$.
Gehe in Endlosschleife.

3. Fall: $f(x_1, \dots, x_n)$ ist undefiniert.

Dann terminiert die Berechnung des Wertes bereits nicht.

□

14. Primitiv rekursive Funktionen und Loop-Programme

In den nächsten beiden Abschnitten betrachten wir Berechnungsmodelle, die eng verwandt sind mit modernen imperativen und funktionalen Programmiersprachen. Die in diesem Abschnitt eingeführten Modelle sind allerdings *nicht* stark genug, um alle (Turing-)berechenbaren Funktionen zu erfassen. Wir werden im Abschnitt 15 zeigen, wie man diese Modelle erweitern muss, um berechnungsuniverselle Modelle zu erhalten.

Wir betrachten hier nur Funktionen von

$$\mathbb{N}^n \rightarrow \mathbb{N}.$$

Dies entspricht dem Spezialfall $|\Sigma| = 1$ bei Wortfunktionen, ist aber keine echte Einschränkung, da es berechenbare Kodierungsfunktionen gibt, d.h.

$$\pi : \Sigma^* \rightarrow \mathbb{N} \quad \text{bijektiv}$$

mit π und π^{-1} berechenbar.

Definition 14.1 (Grundfunktionen)

Die folgenden Funktionen sind *primitiv rekursive Grundfunktionen*:

- 1) $s : \mathbb{N} \rightarrow \mathbb{N}$ mit $x \mapsto x + 1$ (*Nachfolgerfunktion*)
- 2) Für alle $n \geq 0$ und $i, 1 \leq i \leq n$:
 - $\pi_i^{(n)} : \mathbb{N}^n \rightarrow \mathbb{N}$ mit $(x_1, \dots, x_n) \mapsto x_i$ (*Projektion*)
 - $\text{null}^{(n)} : \mathbb{N}^n \rightarrow \mathbb{N}$ mit $(x_1, \dots, x_n) \mapsto 0$ (*Nullfunktion*)

Aus diesen einfachen Funktionen kann man mit Hilfe von Operatoren komplexere Funktionen aufbauen.

Definition 14.2 (Komposition)

Die Funktion $f : \mathbb{N}^n \rightarrow \mathbb{N}$ entsteht aus

$$g : \mathbb{N}^m \rightarrow \mathbb{N} \quad \text{und} \\ h_1, \dots, h_m : \mathbb{N}^n \rightarrow \mathbb{N}$$

durch *Komposition*, falls für alle $(x_1, \dots, x_n) \in \mathbb{N}^n$ gilt:

$$f(x_1, \dots, x_n) = g(h_1(x_1, \dots, x_n), \dots, h_m(x_1, \dots, x_n)).$$

Wendet man Komposition auf echt partielle Funktionen an, so gilt:

$f(x_1, \dots, x_n)$ ist *undefiniert* gdw.

- eines der $h_i(x_1, \dots, x_n)$ ist *undefiniert* oder
- alle h_i -Werte sind *definiert*, aber g von diesen Werten ist *undefiniert*.

Beispiel:

$g(x, y) = x$, $h_1(0) = 0$, $h_2(0)$ undefiniert.

Dann ist $g(h_1(0), h_2(0))$ undefiniert (entspricht call-by-value-Auswertung).

Definition 14.3 (primitive Rekursion)

Es sei $n \geq 0$. Die Funktion $f : \mathbb{N}^{n+1} \rightarrow \mathbb{N}$ entsteht aus

$$\begin{aligned} g : \mathbb{N}^n &\rightarrow \mathbb{N} \text{ und} \\ h : \mathbb{N}^{n+2} &\rightarrow \mathbb{N} \end{aligned}$$

durch *primitive Rekursion*, falls gilt:

- $f(x_1, \dots, x_n, 0) = g(x_1, \dots, x_n)$
- $f(x_1, \dots, x_n, y + 1) = h(x_1, \dots, x_n, f(x_1, \dots, x_n, y), y)$

Auch hier setzen sich wie bei der Komposition wieder undefinierte Werte fort.

Durch Induktion über die letzte Komponente zeigt man leicht, dass dieses Schema die Funktion f eindeutig definiert.

Beispiel 14.4 (Addition)

Die Addition natürlicher Zahlen kann durch primitive Rekursion wie folgt definiert werden:

$$\begin{aligned} \text{add}(x, 0) &= x & &= g(x) \\ \text{add}(x, y + 1) &= \text{add}(x, y) + 1 & &= h(x, \text{add}(x, y), y) \end{aligned}$$

Das heißt also: add entsteht durch primitive Rekursion aus den Funktionen

- $g : \mathbb{N} \rightarrow \mathbb{N}$ mit $x \mapsto x$,
d.h. $g = \pi_1^{(1)}$ ist Grundfunktion,
- $h : \mathbb{N}^3 \rightarrow \mathbb{N}$ mit $(x, z, y) \mapsto z + 1$,
d.h. $h(x, z, y) = s(\pi_2^{(3)}(x, z, y))$.
Also entsteht h durch Komposition aus Grundfunktionen.

Definition 14.5 (Klasse der primitiv rekursiven Funktionen)

Die *Klasse der primitiv rekursiven Funktionen* besteht aus allen Funktionen, die man aus den *Grundfunktionen* durch endlich oftigen Anwenden von

- *Komposition* und
- *primitiver Rekursion*

erhält.

Offenbar sind die Grundfunktionen *total* und die Operationen Komposition und primitive Rekursion erzeugen aus totalen Funktionen wieder totale.

Satz 14.6

Die Klasse der primitiv rekursiven Funktionen enthält nur totale Funktionen.

Trotzdem sind aber die Operationen primitive Rekursion und Komposition auch für partielle Funktionen definiert, und wir werden sie später auch auf partielle Funktionen anwenden.

Beispiel 14.4 zeigt, dass `add` zur Klasse der primitiv rekursiven Funktionen gehört. Wir betrachten nun weitere Beispiele für primitiv rekursive Funktionen.

Multiplikation erhält man durch primitive Rekursion aus der Addition:

$$\begin{aligned} \text{mult}(x, 0) &= 0 && = \text{null}^{(1)}(x) \\ \text{mult}(x, y + 1) &= \text{add}(x, \text{mult}(x, y)) && = \text{add}(\pi_1^{(3)}, \pi_2^{(3)})(x, \text{mult}(x, y), y) \end{aligned}$$

Exponentiation erhält man durch primitive Rekursion aus der Multiplikation:

$$\begin{aligned} \text{exp}(x, 0) &= 1 && = \text{s}(\text{null}^{(1)}(x)) \\ \text{exp}(x, y + 1) &= \text{mult}(x, \text{exp}(x, y)) \end{aligned}$$

Die Funktion

$$\text{min1} : \mathbb{N} \rightarrow \mathbb{N} \text{ mit } x \mapsto x \dot{-} 1 := \begin{cases} x - 1 & x > 0 \\ 0 & x = 0 \end{cases}$$

ist ebenfalls primitiv rekursiv:

$$\begin{aligned} \text{min1}(0) &= 0 && = \text{null}^{(0)}(), \\ \text{min1}(y + 1) &= y && = \pi_2^{(2)}(\text{min1}(y), y) \end{aligned}$$

Übung:

Zeige, dass

$$\text{min} : \mathbb{N}^2 \rightarrow \mathbb{N} \text{ mit } (x, y) \mapsto x \dot{-} y := \begin{cases} x - y & x \geq y \\ 0 & \text{sonst} \end{cases}$$

primitiv rekursiv ist.

Beispiel 14.7

Aus den bisher betrachteten Funktionen erhält man damit die Funktion

$$c : \mathbb{N}^2 \rightarrow \mathbb{N} \text{ mit } (x, y) \mapsto 2^x \cdot (2y + 1) \dot{-} 1$$

durch Komposition, d.h. `c` ist primitiv rekursiv. Diese Funktion ist interessant, da sie eine Bijektion von $\mathbb{N}^2 \rightarrow \mathbb{N}$ ist, d.h. man kann mit ihr Tupel natürlicher Zahlen in eine natürliche Zahl kodieren.

Lemma 14.8

Die Funktion `c` aus Beispiel 14.7 ist eine Bijektion.

Beweis.

Surjektivität: Es sei $z \in \mathbb{N}$. Dann betrachten wir die größte Zweierpotenz 2^x , die $z + 1$ teilt. Offenbar ist dann $\frac{z+1}{2^x}$ eine ungerade Zahl, d.h. es gibt ein y mit

$$\frac{z+1}{2^x} = 2y + 1.$$

Damit ist $z = 2^x \cdot (2y + 1) - 1 = c(x, y)$.

Injektivität: Offenbar ist die größte Zweierpotenz, die $z + 1$ teilt, eindeutig, d.h. x ist eindeutig durch z bestimmt. Damit ist aber auch y eindeutig bestimmt. \square

Da c eine Bijektion ist, gibt es die Umkehrfunktionen c_0 und c_1 mit der Eigenschaft:

- $c_0(c(x, y)) = x$ und $c_1(c(x, y)) = y$
- $c(c_0(z), c_1(z)) = z$

Diese ergeben sich im Prinzip aus dem Beweis von Lemma 14.8, d.h.

$$c_0(z) = \max\{x \in \mathbb{N} \mid 2^x \mid (z + 1)\}$$

$$c_1(z) = ((\frac{z+1}{2^{c_0(z)}}) - 1) / 2$$

Um zu zeigen, dass c_0, c_1 ebenfalls primitiv rekursiv sind, benötigen wir noch etwas Vorarbeit.

Die Funktionen $\text{sign} : \mathbb{N} \rightarrow \mathbb{N}$ und $\overline{\text{sign}} : \mathbb{N} \rightarrow \mathbb{N}$ mit

$$\text{sign}(x) = \begin{cases} 0 & x = 0 \\ 1 & x > 0 \end{cases} \quad \overline{\text{sign}}(x) = \begin{cases} 1 & x = 0 \\ 0 & x > 0 \end{cases}$$

sind primitiv rekursiv:

$$\begin{array}{ll} \text{sign}(0) = 0 & \text{sign}(y + 1) = 1 \\ \overline{\text{sign}}(0) = 1 & \overline{\text{sign}}(1) = 0 \end{array}$$

sind primitiv rekursive Definitionsschemata dafür.

Definition 14.9 (Fallunterscheidung)

Es seien $g_1, g_2, h : \mathbb{N}^n \rightarrow \mathbb{N}$ gegeben.

Die Funktion $f : \mathbb{N}^n \rightarrow \mathbb{N}$ entsteht daraus durch *Fallunterscheidung*, falls für alle $\underline{x} \in \mathbb{N}^n$ gilt:

$$f(\underline{x}) = \begin{cases} g_1(\underline{x}) & \text{falls } h(\underline{x}) = 0 \\ g_2(\underline{x}) & \text{falls } h(\underline{x}) > 0 \end{cases}$$

Lemma 14.10

Sind g_1, g_2, h primitiv rekursiv, so auch f .

Beweis.

$$f(\underline{x}) = g_1(\underline{x}) \cdot \overline{\text{sign}}(h(\underline{x})) + g_2(\underline{x}) \cdot \text{sign}(h(\underline{x})). \quad \square$$

Definition 14.11 (beschränkte Minimalisierung)

Es sei $n \geq 0$. Die Funktion $f : \mathbb{N}^{n+1} \rightarrow \mathbb{N}$ entsteht aus $g : \mathbb{N}^{n+1} \rightarrow \mathbb{N}$ durch *beschränkte Minimalisierung*, falls gilt:

$$f(\underline{x}, y) = \begin{cases} j & \text{falls } j = \min\{i \leq y \mid g(\underline{x}, i) = 0\} \text{ existiert} \\ y + 1 & \text{sonst} \end{cases}$$

Wir schreiben dann $f = \bar{\mu}g$.

Lemma 14.12

Ist g primitiv rekursiv, so auch $\bar{\mu}g$.

Beweis. Wir definieren $\bar{\mu}g$ mittels primitiver Rekursion wie folgt:

$$1) \quad \bar{\mu}g(\underline{x}, 0) = \begin{cases} 0 & \text{falls } g(\underline{x}, 0) = 0 \\ 1 & \text{sonst} \end{cases}$$

D.h. $\bar{\mu}g(\underline{x}, 0) = \text{sign}(g(\underline{x}, 0))$

$$2) \quad \bar{\mu}g(\underline{x}, y + 1) = \begin{cases} \bar{\mu}g(\underline{x}, y) & \text{falls } \bar{\mu}g(\underline{x}, y) \leq y \text{ oder } g(\underline{x}, y + 1) = 0 \\ y + 2 & \text{sonst} \end{cases}$$

Es handelt sich hier also um eine Fallunterscheidung. Es bleibt zu zeigen, dass die Funktion

$$h(\underline{x}, z, y) = \begin{cases} 0 & \text{falls } z \leq y \text{ oder } g(\underline{x}, y + 1) = 0 \\ 1 & \text{sonst} \end{cases}$$

primitiv rekursiv ist.

Offenbar ist aber

$$h(\underline{x}, z, y) = \text{sign}((z \dot{-} y) \cdot g(\underline{x}, y + 1))$$

und damit primitiv rekursiv. □

Wir kommen nun zurück zu den Umkehrfunktionen der Bijektion aus Beispiel 14.7. Betrachten wir zunächst die Teilbarkeitsrelation in der Definition von c_0 .

Lemma 14.13

Die Funktion

$$\text{teilt} : \mathbb{N}^2 \rightarrow \mathbb{N} \text{ mit } (x, y) \mapsto \begin{cases} 0 & x|y \\ 1 & \text{sonst} \end{cases}$$

ist primitiv rekursiv.

Beweis. Die Funktion

$$h(x, z, y) = \begin{cases} 0 & x \cdot z = y \\ 1 & \text{sonst} \end{cases}$$

ist primitiv rekursiv, da

$$h(x, y, z) = \text{sign}((x \cdot z \dot{-} y) + (y \dot{-} x \cdot z)).$$

Damit ist auch $\bar{\mu}h$ primitiv rekursiv, und es gilt:

$$\bar{\mu}h(x, y, y) = \begin{cases} j & \text{falls } j = \min\{i \leq y \mid x \cdot i = y\} \\ y + 1 & \text{sonst} \end{cases}$$

Damit ist aber $\text{teilt}(x, y) = \bar{\mu}h(x, y, y) \dot{-} y$. □

Lemma 14.14

Die Umkehrfunktion c_0 der Funktion c aus Beispiel 14.7 ist primitiv rekursiv.

Beweis. $c_0(z) = \max\{x \in \mathbb{N} \mid 2^x \mid (z + 1)\}$.

Es genügt dafür, nach dem kleinsten $i \leq z + 1$ zu suchen mit

$$2^{(z+1) \dot{-} i} \mid (z + 1)$$

d.h. nach dem kleinsten $i \leq z + 1$ mit

$$\text{teilt}(2^{(z+1) \dot{-} i}, z + 1) = 0.$$

Dies gelingt durch beschränkte Minimalisierung. □

Um zu zeigen, dass auch c_1 primitiv rekursiv ist, betrachten wir die ganzzahlige Division

$$\text{div} : \mathbb{N}^2 \rightarrow \mathbb{N} \text{ mit } (x, y) \mapsto \begin{cases} \lfloor x/y \rfloor & \text{falls } y > 0 \\ x & \text{sonst} \end{cases}$$

wobei $\lfloor x/y \rfloor$ die größte natürliche Zahl unterhalb von x/y ist.

Beachte:

Ist $y > 0$ und x durch y teilbar, so ist $\text{div}(x, y) = x/y$. In der Definition von c_1 sind diese Bedingungen erfüllt. Es ist daher

$$c_2(z) = \text{div}(\text{div}(z + 1, 2^{c_1(z)}) \dot{-} 1, 2).$$

Lemma 14.15

Die Umkehrfunktion c_1 der Funktion c aus Beispiel 14.7 ist primitiv rekursiv.

Beweis. Es genügt zu zeigen, dass div primitiv rekursiv ist. Wir betrachten dazu die Funktion

$$f(x, y, z) = \begin{cases} 0 & \text{falls } z \cdot y > x \\ 1 & \text{falls } z \cdot y \leq x \end{cases}$$

Diese ist primitiv rekursiv, da

$$f(x, y, z) = \overline{\text{sign}}(z \cdot y \dot{-} x).$$

Damit ist auch $\bar{\mu}f$ primitiv rekursiv und somit (unter Verwendung geeigneter Projektionen) auch

$$g(x, y) := \bar{\mu}f(x, y, x).$$

Es gilt nun aber

$$g(x, y) = \begin{cases} \text{das kleinste } i \leq x \text{ mit } i \cdot y > x & \text{falls existent} \\ x + 1 & \text{sonst} \end{cases}$$

Ist $y > 1$, so existiert so ein i stets und es ist $\text{div}(x, y) = i - 1$.

Ist $y = 1$ oder $y = 0$, so existiert so ein i nicht, d.h. $x + 1$ wird ausgegeben. In diesen Fällen ist aber auch $\text{div}(x, y) = x$. Also gilt:

$$\text{div}(x, y) = g(x, y) - 1. \quad \square$$

Wir betrachten nun eine einfache imperative Programmiersprache, die genau die primitiv rekursiven Funktionen berechnen kann.

LOOP-Programme sind aus den folgenden Komponenten aufgebaut:

- Variablen: x_0, x_1, x_2, \dots
- Konstanten: $0, 1, 2, \dots$ (also die Elemente von \mathbb{N})
- Trennsymbole: ; und :=
- Operationssymbole: + und -
- Schlüsselwörter: LOOP, DO, END

Definition 14.16 (Syntax LOOP)

Die *Syntax von LOOP-Programmen* ist induktiv definiert:

- 1) Jede *Wertzuweisung*

$$\begin{aligned} x_i &:= x_j + c \text{ und} \\ x_i &:= x_j - c \end{aligned}$$

für $i, j \geq 0$ und $c \in \mathbb{N}$ ist ein LOOP-Programm.

- 2) Falls P_1 und P_2 LOOP-Programme sind, so ist auch

$$P_1; P_2 \quad (\text{Hintereinanderausführung})$$

ein LOOP-Programm.

- 3) Falls P ein LOOP-Programm ist und $i \geq 0$, so ist auch

$$\text{LOOP } x_i \text{ DO } P \text{ END}$$

ein LOOP-Programm.

Die *Semantik* dieser einfachen Sprache ist wie folgt definiert:

Bei einem LOOP-Programm, das eine Funktion $f : \mathbb{N}^k \rightarrow \mathbb{N}$ berechnen soll:

- werden die *Variablen* x_1, \dots, x_k mit den *Eingabewerten* n_1, \dots, n_k vorbesetzt.
- Alle anderen Variablen erhalten den Wert 0.
- *Ausgabe* ist der Wert der Variablen x_0 nach Ausführung des Programms.

Die einzelnen Programmkonstrukte haben die folgende Bedeutung:

- 1) $x_i := x_j + c$
 Der neue Wert der Variablen x_i ist die Summe des alten Wertes von x_j und c .
 $x_i := x_j - c$
 Der neue Wert der Variablen x_i ist der Wert von x_j minus c , falls dieser Wert ≥ 0 ist und 0 sonst.
- 2) $P_1; P_2$
 Hier wird zunächst P_1 und dann P_2 ausgeführt.
- 3) LOOP x_i DO P END
 Das Programm P wird sooft ausgeführt, wie der Wert von x_i zu Beginn angibt. Änderungen des Wertes von x_i während der Ausführung von P haben keinen Einfluss auf die Anzahl der Schleifendurchläufe.

Definition 14.17 (LOOP-berechenbar)

Die Funktion

$$f : \mathbb{N}^k \rightarrow \mathbb{N}$$

heißt *LOOP-berechenbar*, falls es ein LOOP-Programm P gibt, das f in dem folgenden Sinne berechnet:

- *Gestartet* mit n_1, \dots, n_k in den Variablen x_1, \dots, x_k (und 0 in den restlichen Variablen)
- *stoppt* P mit dem Wert $f(n_1, \dots, n_k)$ in der Variablen x_0 .

Bemerkung:

Offenbar terminieren LOOP-Programme stets, da für jede Schleife eine feste Anzahl von Durchläufen durch den anfänglichen Wert der Schleifenvariablen festgelegt wird. Daher sind alle durch LOOP-Programme berechneten Funktionen total.

Beispiel:

Die Additionsfunktion ist LOOP-berechenbar:

```

 $x_0 := x_1 + 0;$ 
LOOP  $x_2$  DO  $x_0 := x_0 + 1$  END
    
```

Mit den Programmkonstrukten von Definition 14.16 kann man auch andere in Programmiersprachen vorhandene Konstrukte simulieren:

```

IF  $x = 0$  THEN  $P$  END
    
```

kann z.B. simuliert werden durch

```

 $y := 1;$ 
LOOP  $x$  DO  $y := 0$  END;
LOOP  $y$  DO  $P$  END
    
```

wobei y eine neue Variable ist, die nicht in P vorkommt und $\neq x_0$ ist.

Satz 14.18

Die Klasse der primitiv rekursiven Funktionen stimmt mit der der LOOP-berechenbaren Funktionen überein.

Beweis.

(I) Alle primitiv rekursiven Funktionen sind LOOP-berechenbar:

- Für die *Grundfunktionen* ist klar, dass sie LOOP-berechenbar sind.
- Komposition:

$$f(x_1, \dots, x_n) = g(h_1(x_1, \dots, x_k), \dots, h_m(x_1, \dots, x_k))$$

Es seien P_1, \dots, P_m, P LOOP-Programme für h_1, \dots, h_m, g .

Durch Speichern der Eingabewerte und der Zwischenergebnisse in unbenutzten Variablen kann man zunächst die Werte von h_1, \dots, h_m mittels P_1, \dots, P_m berechnen und dann auf diese Werte P anwenden.

- primitive Rekursion:

$$\begin{aligned} f(x_1, \dots, x_n, 0) &= g(x_1, \dots, x_n) \\ f(x_1, \dots, x_n, x_{n+1} + 1) &= h(x_1, \dots, x_n, f(x_1, \dots, x_n, x_{n+1}), x_{n+1}) \end{aligned}$$

Die Funktion f kann durch das LOOP-Programm

```

z1 := g(x1, ..., xn);   (★)
z2 := 0;
z3 := xn+1 - 1;
LOOP z3 DO
  z1 := h(x1, ..., xn, z1, z2);   (★)
  z2 := z2 + 1
END
    
```

berechnet werden.

Dabei sind die mit (★) gekennzeichneten Anweisungen Abkürzungen für Programme, welche Ein- und Ausgaben geeignet kopieren und die Programme für g und h anwenden.

Die Variablen z_1, z_2, z_3 sind neue Variablen, die in den Programmen für g und h *nicht* vorkommen.

(II) Alle LOOP-berechenbaren Funktionen sind primitiv rekursiv.

Dazu beschaffen wir uns zunächst eine primitiv rekursive Bijektion

$$c^{(n)} : \mathbb{N}^n \rightarrow \mathbb{N} \quad (n \geq 2)$$

sowie die zugehörigen Umkehrfunktionen

$$c_0^{(n)}, \dots, c_{n-1}^{(n)}$$

Wir definieren dazu

$$\begin{aligned}
 c^{(n)}(x_1, \dots, x_n) &:= c(x_1, c(x_2, \dots, c(x_{n-1}, x_n) \dots)) \\
 c_0^{(n)}(z) &:= c_0(z) \\
 c_1^{(n)}(z) &:= c_1(c_2(z)) \\
 &\vdots \\
 c_{n-2}^{(n)}(z) &:= c_1(c_2^{(n-2)}(z)) \\
 c_{n-1}^{(n)}(z) &:= c_2^{n-1}(z)
 \end{aligned}$$

Da c und c_1, c_2 primitiv rekursiv sind, sind auch $c^{(n)}$ und $c_0^{(n)}, \dots, c_{n-1}^{(n)}$ primitiv rekursiv.

Es sei nun P ein LOOP-Programm, das die Funktion $f : \mathbb{N}^r \rightarrow \mathbb{N}$ berechnet. Es sei l der maximale Index der in P vorkommenden Variablen und $k := \max\{r, l\}$.

Wir zeigen durch Induktion über den Aufbau von LOOP-Programmen, dass die Funktion $g_P : \mathbb{N} \rightarrow \mathbb{N}$ primitiv rekursiv ist, wobei

$$g_P(z) = c^{(k+1)}(b_0, \dots, b_k)$$

wenn

- b_0, \dots, b_k die Werte von x_0, \dots, x_k nach Ausführung des Programms P
- bei Startwerten $a_0 = c_0^{(k+1)}(z), \dots, a_k = c_k^{(k+1)}(z)$ sind.

1) Hat P die Form $x_i := x_j + c$, so ist

$$g_P(z) = c^{(k+1)}(c_0^{(k+1)}(z), \dots, c_{i-1}^{(k+1)}(z), c_j^{(k+1)}(z) + c, c_{i+1}^{(k+1)}(z), \dots).$$

Primitiv rekursiv.

Entsprechend kann die andere Zuweisung behandelt werden.

2) $P = Q; R$

Dann ist

$$g_P(z) = g_R(g_Q(z)).$$

Da

- g_Q, g_R nach Induktionsvoraussetzung primitiv rekursiv sind und
- g_P durch Komposition daraus entsteht

ist auch g_P primitiv rekursiv.

3) $P = \text{LOOP } x_i \text{ DO } Q \text{ END}$

Wir definieren zunächst die zweistellige Funktion h durch primitive Rekursion aus g_Q :

$$\begin{aligned}
 h(x, 0) &= x, \\
 h(x, y + 1) &= g_Q(h(x, y))
 \end{aligned}$$

Offenbar liefert $h(z, n)$

- die Kodierung der Werte der Variablen x_0, \dots, x_k ,
- nachdem man, beginnend mit $x_0 = c_0^{(k+1)}(z), \dots, x_k = c_k^{(k+1)}(z)$,
- das Programm Q n -mal ausgeführt hat.

Daher gilt:

$$g_P(z) = h(z, c_i^{(k+1)}(z)).$$

Dies schließt den Induktionsbeweis ab, dass die Funktion g_P für jedes LOOP-Programm P primitiv rekursiv ist.

Ist

$$f : \mathbb{N}^r \rightarrow \mathbb{N}$$

die von P berechnete Funktion, so gilt:

$$f(z_1, \dots, z_r) = c_0^{(k+1)}(g_P(c^{(k+1)}(0, z_1, \dots, z_r, \underbrace{0, \dots, 0}_{k-r}))). \quad \square$$

Da wir durch LOOP-berechenbare/primitiv rekursive Funktionen nur totale Funktionen erhalten, ist *nicht* jede (intuitiv) berechenbare Funktion in dieser Klasse enthalten. Aber was ist mit den totalen berechenbaren Funktionen?

Satz 14.19

Es gibt totale berechenbare Funktionen, die nicht LOOP-berechenbar sind.

Beweis. Wir definieren die Länge von LOOP-Programmen induktiv über deren Aufbau:

- 1) $|x_i := x_j + c| := i + j + c + 1$
 $|x_i := x_j - c| := i + j + c + 1$
- 2) $|P; Q| := |P| + |Q| + 1$
- 3) $|\text{LOOP } x_i \text{ DO } P \text{ END}| := |P| + i + 1$

Für eine gegebene Länge n gibt es nur endlich viele LOOP-Programme dieser Länge. Deshalb macht die folgende Definition Sinn:

$$f(x, y) := 1 + \max\{g(y) \mid g : \mathbb{N} \rightarrow \mathbb{N} \text{ wird von einem LOOP-Programm der Länge } x \text{ berechnet}\}$$

Behauptung 1:

Die Funktion

$$d : \mathbb{N} \rightarrow \mathbb{N} \text{ mit } z \mapsto f(z, z)$$

ist nicht LOOP-berechenbar, denn:

Sei P ein LOOP-Programm, das d berechnet und sei $n = |P|$.

Wir betrachten $d(n) = f(n, n)$. Nach Definition ist $f(n, n)$ größer als der maximale Funktionswert, den ein LOOP-Programm der Länge n bei Eingabe n berechnen kann. Dies widerspricht der Tatsache, dass d von einem LOOP-Programm der Länge n berechnet wird.

Behauptung 2:

Die Funktion d ist (intuitiv) berechenbar, denn:

Bei Eingabe z zählt man die *endlich vielen* LOOP-Programme der Länge z auf und wendet sie jeweils auf die Eingabe z an. Da alle diese Aufrufe terminieren, kann man in endlicher Zeit den maximalen so erhaltenen Funktionswert berechnen.

Mit der Churchschen These ist diese Funktion auch Turing-berechenbar. □

Eine prominente *nicht* primitiv rekursive, aber berechenbare Funktion ist die sogenannte *Ackermannfunktion* A , die wie folgt definiert ist:

$$\begin{aligned} A(0, y) &:= y + 1 \\ A(x + 1, 0) &:= A(x, 1) \\ A(x + 1, y + 1) &:= A(x, A(x + 1, y)) \end{aligned}$$

Durch (geschachtelte) Induktion zeigt man leicht, dass A dadurch eindeutig definiert ist. Man kann zeigen, dass A zwar Turing-berechenbar ist, aber nicht LOOP-berechenbar, da sie schneller wächst als jede LOOP-berechenbare Funktion (siehe [Schö97]).

15. μ -rekursive Funktionen und While-Programme

Um alle berechenbaren Funktionen zu erhalten, muss man die primitiv rekursiven Funktionen um eine weitere Operation, die *unbeschränkte Minimalisierung*, erweitern. Die LOOP-Programme muss man um eine *While-Schleife* (d.h. eine Schleife ohne vorher festgelegte Anzahl der Durchläufe) erweitern.

Definition 15.1 (unbeschränkte Minimalisierung)

Es sei $n \geq 0$. Die Funktion $f : \mathbb{N}^n \rightarrow \mathbb{N}$ entsteht aus $g : \mathbb{N}^{n+1} \rightarrow \mathbb{N}$ durch *unbeschränkte Minimalisierung* (Anwendung des μ -Operators), falls gilt:

$$f(x_1, \dots, x_n) = \begin{cases} y & \text{falls } g(x_1, \dots, x_n, y) = 0 \text{ und} \\ & g(x_1, \dots, x_n, z) \text{ ist definiert und } \neq 0 \\ & \text{für alle } 0 \leq z < y \\ \text{undefiniert} & \text{sonst} \end{cases}$$

Wir schreiben dann auch $f = \mu g$.

Beachte:

Der μ -Operator sucht nach dem kleinsten y , so dass $g(\underline{x}, y) = 0$ ist. Dabei müssen aber alle vorherigen Werte $g(\underline{x}, z)$ für $z < y$ definiert sein.

Beispiel:

$$1) \ g_1(x, y) = \begin{cases} x - y & \text{falls } x \geq y \\ \text{undefiniert} & \text{falls } x < y \end{cases}$$

$$\mu g_1(x) = x,$$

d.h. man erhält durch Anwenden des μ -Operators auf die partielle Funktion g_1 die totale Identitätsfunktion $f_1 : \mathbb{N} \rightarrow \mathbb{N}$ mit $x \mapsto x$.

$$2) \ g_2(x, y) = \begin{cases} y - x & \text{falls } x \leq y \\ \text{undefiniert} & \text{falls } x > y \end{cases}$$

$$\mu g_2(x) = \begin{cases} 0 & \text{falls } x = 0 \\ \text{undefiniert} & \text{falls } x > 0 \end{cases} \quad (\text{vor Wert 0 hat man undefiniert})$$

$$3) \ g_3(x, y) = x + y$$

$$\mu g_3(x) = \begin{cases} 0 & \text{falls } x = 0 \\ \text{undefiniert} & \text{falls } x > 0 \end{cases} \quad (\text{kein Wert 0})$$

Durch Anwenden von μ kann man also auch partielle Funktionen aus totalen erhalten.

Definition 15.2 (Klasse der μ -rekursiven Funktionen)

Die Klasse der μ -rekursiven Funktionen besteht aus den Funktionen, welche man aus den Grundfunktionen (Definition 14.1) durch endlich oftigen Anwenden von

- Komposition,
- primitiver Rekursion und
- unbeschränkter Minimalisierung

erhält.

Definition 15.3 (Syntax von WHILE-Programmen)

Die Syntax von WHILE-Programmen enthält alle Konstrukte in der Syntax von LOOP-Programmen und zusätzlich

- 4) Falls P ein WHILE-Programm ist und $i \geq 0$, so ist auch
- $$\text{WHILE } x_i \neq 0 \text{ DO } P \text{ END}$$
- ein WHILE-Programm.

Die Semantik dieses Konstrukts ist wie folgt definiert:

- Das Programm P wird solange iteriert, bis x_i den Wert 0 erhält.
- Geschieht das nicht, so terminiert diese Schleife nicht.

Beachte:

Man könnte bei der Definition der WHILE-Programme auf das LOOP-Konstrukt verzichten, da es durch WHILE simulierbar ist:

$$\text{LOOP } x \text{ DO } P \text{ END}$$

kann simuliert werden durch:

$$y := x + 0;$$
$$\text{WHILE } y \neq 0 \text{ DO } y := y - 1; P \text{ END}$$

wobei y eine neue Variable ist.

Definition 15.4 (WHILE-berechenbar)

Die (partielle) Funktion $f : \mathbb{N}^k \rightarrow \mathbb{N}$ heißt *WHILE-berechenbar*, falls es ein WHILE-Programm P gibt, das f in dem folgenden Sinne berechnet:

- *Gestartet* mit n_1, \dots, n_k in den Variablen x_1, \dots, x_k (und 0 in den restlichen Variablen)
- *stoppt* P mit dem Wert $f(n_1, \dots, n_k)$ in der Variablen x_0 , falls dieser Wert definiert ist.
- Sonst stoppt P nicht.

Beispiel:

Die Funktion

$$f : \mathbb{N}^2 \rightarrow \mathbb{N} \text{ mit } (x, y) \mapsto \begin{cases} x - y & \text{falls } x \geq y \\ \text{undefiniert} & \text{sonst} \end{cases}$$

ist WHILE-berechenbar durch das folgende Programm:

```
WHILE  $x_2 \neq 0$  DO
   $y := x_1$ ;
  WHILE  $y \neq 0$  DO
     $x_1 := x_1 - 1$ ;
     $x_2 := x_2 - 1$ ;
     $y := 0$ 
  END
END;
 $x_0 := x_1$ 
```

Satz 15.5

Die Klasse der μ -rekursiven Funktionen stimmt genau mit der Klasse der WHILE-berechenbaren Funktionen überein.

Beweis. Wir müssen hierzu den Beweis von Satz 14.18 um die Behandlung des zusätzlichen Operators/Konstrukts ergänzen.

- 1) Alle μ -rekursiven Funktionen sind WHILE-berechenbar:

Es sei $f = \mu g$ und P (nach Induktionsvoraussetzung) ein WHILE-Programm für g . Dann berechnet das folgende Programm die Funktion f :

```
 $x_0 := 0$ ;
 $y := g(x_1, \dots, x_n, x_0)$ ; (Realisierbar mittels  $P$ )
WHILE  $y \neq 0$  DO
   $x_0 := x_0 + 1$ ;
   $y := g(x_1, \dots, x_n, x_0)$ 
END
```

Beachte:

Dieses Programm ist nur deshalb korrekt, weil wir in der Definition von μ gefordert haben, dass alle vorherigen Werte nicht nur $\neq 0$, sondern auch definiert sein müssen.

- 2) Alle WHILE-berechenbaren Funktionen sind μ -rekursiv:

Betrachte das Programm `WHILE $x_i \neq 0$ DO Q END.`

Wie im Beweis von Satz 14.18 bei der Behandlung von LOOP beschaffen wir uns eine μ -rekursive Funktion

$$h : \mathbb{N}^2 \rightarrow \mathbb{N},$$

für die

- $h(z, n)$ die Kodierung der Werte der Variablen x_0, \dots, x_k ist, nachdem man,
- beginnend mit $x_0 = c_0^{(k+1)}(z), \dots, x_k = c_k^{(k+1)}(z)$
- das Programm Q n -mal ausgeführt hat.

Der μ -Operator, angewandt auf die Komposition $c_i^{(k+1)}(h)$, sucht nach der kleinsten Iterationszahl, so dass die Variable $x_i = 0$ wird. Daher ist

$$g_P(z) = h(z, (\mu c_i^{(k+1)}(h(z))))). \quad \square$$

Es bleibt noch zu zeigen, dass die Klasse der μ -rekursiven/WHILE-berechenbaren Funktionen mit der Klasse der Turing-berechenbaren Funktionen übereinstimmt.

Satz 15.6

Jede μ -rekursive Funktion ist Turing-berechenbar.

Beweis. Es ist leicht zu zeigen, dass die Grundfunktionen Turing-berechenbar sind. (Übung)

Komposition: Seien $\mathcal{A}_g, \mathcal{A}_{h_1}, \dots, \mathcal{A}_{h_m}$ DTM, die g, h_1, \dots, h_m berechnen.

Verwende $m + 1$ Bänder:

- Kopiere die Eingabe \underline{x} auf Band2, ..., Band($m + 1$)
- \mathcal{A}_{h_i} berechnet $h_i(\underline{x})$ auf Band ($i + 1$)
- Überschreibe Band1 mit $\# a^{h_1(\underline{x})} \# \dots \# a^{h_m(\underline{x})} \#$
- Berechne g -Wert davon mit \mathcal{A}_g

Primitive Rekursion: Seien $\mathcal{A}_g, \mathcal{A}_h$ DTM für g, h .

Verwende 4 Bänder:

- Band1: Speichert \underline{xy}
- Band2: Zählt von 0 ab hoch bis y (aktueller Wert: z)
- Band3: Enthält $f(\underline{x}, z)$ (für $z = 0$ mittels \mathcal{A}_g berechnet)
- Band4: Berechnung von \mathcal{A}_h

μ -Operator: Sei \mathcal{A}_g eine DTM für g .

Verwende 3 Bänder:

- Band1: Speichert Eingabe \underline{x}
- Band2: Zählt von 0 ab hoch (aktueller Wert: z)
- Band3: Berechne $g(\underline{x}, z)$ mittels \mathcal{A}_g und teste, ob Wert = 0 ist

□

Satz 15.7

Jede Turing-berechenbare Funktion ist WHILE-berechenbar.

Beweisskizze. Um diesen Satz zu beweisen, müssen wir Konfigurationen von Turingmaschinen der Form

$$\alpha q \beta \text{ für } \alpha, \beta \in \Gamma^+ \text{ und } q \in Q$$

in drei natürlichen Zahlen kodieren.

Diese werden dann in den drei Programmvariablen x_1, x_2, x_3 des WHILE-Programms gespeichert:

- x_1 repräsentiert α ,
- x_2 repräsentiert q ,
- x_3 repräsentiert β .

Es sei o.B.d.A. $\Gamma = \{a_1, \dots, a_n\}$ und $Q = \{q_1, \dots, q_k\}$.

Die Konfiguration

$$a_{i_1} \dots a_{i_l} q_m a_{j_1} \dots a_{j_r}$$

wird dargestellt durch

$$\begin{aligned} x_1 &= (i_1, \dots, i_l)_b & := \sum_{\nu=1}^l i_\nu \cdot b^{l-\nu} \\ x_2 &= m \\ x_3 &= (j_r, \dots, j_1)_b & := \sum_{\rho=1}^r j_\rho \cdot b^{\rho-1}, \end{aligned}$$

wobei $b > |\Gamma|$ ist, d.h.

- $a_{i_1} \dots a_{i_l}$ repräsentiert x_1 in b -närer Zahlendarstellung und
- $a_{j_r} \dots a_{j_1}$ (Reihenfolge!) repräsentiert x_3 in b -närer Zahlendarstellung.

Herauslesen des aktuellen Symbols a_{j_1}

Ist $x_3 = (j_r, \dots, j_1)_b$, so ist $j_1 = x_3 \bmod b$.

Ändern dieses Symbols zu a_j

Der neue Wert von x_3 ist $(j_r, \dots, j_2, j)_b = \text{div}((j_r, \dots, j_2, j_1)_b, b) \cdot b + j$

Verschieben des Schreib-Lesekopfes

kann durch ähnliche arithmetische Operationen realisiert werden.

All diese Operationen sind offensichtlich WHILE-berechenbar (sogar LOOP!).

Das WHILE-Programm, welches die gegebene DTM simuliert, arbeitet wie folgt:

- 1) Aus der Eingabe wird die Kodierung der Startkonfiguration der DTM in den Variablen x_1, x_2, x_3 erzeugt.
- 2) In einer WHILE-Schleife wird bei jedem Durchlauf ein Schritt der TM-Berechnung simuliert (wie oben angedeutet)
 - In Abhängigkeit vom aktuellen Zustand (Wert von x_2) und
 - dem gelesenen Symbol, d.h. von $x_3 \bmod b$
 - wird das aktuelle Symbol verändert und
 - der Schreib-Lesekopf bewegt.

Die WHILE-Schleife terminiert, wenn der aktuelle Zustand zusammen mit dem gelesenen Symbol keinen Nachfolgezustand hat.

All dies ist durch einfache (WHILE-berechenbare) arithmetische Operationen möglich.

- 3) Aus dem Wert von x_3 nach Termination der WHILE-Schleife wird der Ausgabewert herausgelesen und in die Variable x_0 geschrieben. □

Insgesamt haben wir also gezeigt:

Theorem 15.8

Die folgenden Klassen von Funktionen stimmen überein:

- 1) Turing-berechenbare Funktionen*
- 2) WHILE-berechenbare Funktionen*
- 3) μ -rekursive Funktionen*

Diese Äquivalenz ist ein wichtiges Argument für die Korrektheit der Churchschen These.

16. Universelle Maschinen und unentscheidbare Probleme

Wir werden hier zeigen, dass es Relationen gibt, die nicht Turing-entscheidbar sind, d.h. Relationen, deren charakteristische Funktion nicht Turing-berechenbar ist. Mit der Churchschen These sind diese Relationen dann nicht entscheidbar im intuitiven Sinn.

Dazu konstruieren wir eine *universelle Turingmaschine*, die alle Turingmaschinen simulieren kann. Die universelle Maschine erhält als zusätzliche Eingabe eine Kodierung der zu simulierenden Turingmaschine.

Konventionen:

- *Arbeitsalphabete* der betrachteten Turingmaschinen sind Teilmengen von $\{a_1, a_2, a_3, \dots\}$, wobei $a_1 = a$, $a_2 = b$, $a_3 = \beta$.
- *Zustandsmengen* sind Teilmengen von $\{q_1, q_2, q_3, \dots\}$, wobei q_1 stets der *Anfangszustand* ist.

Definition 16.1 (Kodierung einer Turingmaschine)

Es sei $\mathcal{A} = (Q, \Sigma, \Gamma, q_1, \Delta, F)$ eine Turingmaschine, die o.B.d.A. die obigen Konventionen erfüllt.

- 1) Eine *Transition*

$$t = (q_i, a_j, a_{j'}, \overset{l}{r}, q_{i'})$$

wird *kodiert* durch

$$\text{code}(t) = a^i b a^j b a^{j'} b \overset{a}{aa} b a^{i'} b b.$$

aaa

- 2) Besteht Δ aus den Transitionen t_1, \dots, t_k und ist $F = \{q_{i_1}, \dots, q_{i_r}\}$, so wird \mathcal{A} *kodiert* durch

$$\text{code}(\mathcal{A}) = \text{code}(t_1) \dots \text{code}(t_k) b a^{i_1} b \dots b a^{i_r} b b b.$$

Bemerkung 16.2.

- 1) Für $x = \text{code}(\mathcal{A})w$ mit $w \in \Sigma^*$ folgt w genau auf den zweiten Block bbb , kann also aus x eindeutig wieder herausgelesen werden.
- 2) Es gibt eine DTM \mathcal{A}_{CODE} , welche die Relation

$$CODE = \{\text{code}(\mathcal{A}) \mid \mathcal{A} \text{ ist DTM über } \Sigma\}$$

entscheidet, denn:

- Überprüfe bei Eingabe w zunächst, ob w eine Turingmaschine kodiert (d.h. eine Folge von Transitionskodierungen gefolgt von einer Endzustandsmengenkodierung ist).
- Überprüfe dann, ob die kodierte Turingmaschine deterministisch ist (bei jeder Transition wird nachgeschaut, ob es eine andere mit demselben Anfangsteil gibt).

Satz 16.3 (Turing)

Es gibt eine universelle DTM \mathcal{U} über Σ , d.h. eine DTM mit der folgenden Eigenschaft: Für alle DTM \mathcal{A} und alle $w \in \Sigma^*$ gilt:

$$\mathcal{U} \text{ akzeptiert } \text{code}(\mathcal{A})w \text{ gdw. } \mathcal{A} \text{ akzeptiert } w.$$

Es ist also \mathcal{U} ein „Turing-Interpreter für Turingmaschinen“.

Beweis. \mathcal{U} führt bei Eingabe $\text{code}(\mathcal{A})w$ die \mathcal{A} -Berechnung

$$k_0 \vdash_{\mathcal{A}} k_1 \vdash_{\mathcal{A}} \dots$$

in kodierter Form aus, d.h. \mathcal{U} erzeugt sukzessive Bandbeschriftungen

$$\text{code}(\mathcal{A})\text{code}(k_0), \text{code}(\mathcal{A})\text{code}(k_1), \text{code}(\mathcal{A})\text{code}(k_2), \dots$$

Konfigurationskodierung:

$$\text{code}(a_{i_1} \dots a_{i_l} q_j a_{i_{l+1}} \dots a_{i_r}) = a^{i_1} b \dots a^{i_l} b a^j b b a^{i_{l+1}} b \dots a^{i_r} b$$

Arbeitsweise von \mathcal{U} : (bei Eingabe $\text{code}(\mathcal{A})w$)

- Erzeuge aus $\text{code}(\mathcal{A})w$ die Anfangsbeschriftung

$$\text{code}(\mathcal{A})\text{code}(\underbrace{b q_1 w b}_{k_0}).$$

- Simuliere die \mathcal{A} -Schritte, ausgehend vom jeweiligen Konfigurationskode

$$\dots b a^j b b a^i b \dots \text{ (Zustand } q_j, \text{ gelesenes Symbol } a_i).$$

- Aufsuchen einer Transitionskodierung $\text{code}(t)$, die mit $a^j b a^i$ beginnt.
- Falls es so eine gibt, Änderungen der Konfigurationskodierung entsprechend t .
- Sonst geht \mathcal{U} in Stoppzustand. Dieser ist akzeptierend gdw. a^j in der Kodierung der Endzustandsmenge vorkommt.

□

Satz 16.4

Die Relation

$$UNIV = \{\text{code}(\mathcal{A})w \in \Sigma^* \mid \mathcal{A} \text{ ist DTM über } \Sigma, \text{ die } w \text{ akzeptiert}\}$$

ist partiell entscheidbar, aber nicht entscheidbar.

Beweis.

1) Wir zeigen, dass $UNIV$ Turing-akzeptierbar (und damit partiell entscheidbar) ist:

Bei Eingabe x geht die TM, welche $UNIV$ akzeptiert, wie folgt vor:

- Teste, ob x von der Form $x = x_1w$ ist mit

$$x_1 \in \text{CODE und } w \in \Sigma^*.$$

- Wenn ja, so wende \mathcal{U} auf x an.

2) Angenommen, $UNIV$ ist rekursiv. Dann ist auch $\overline{UNIV} = \Sigma^* \setminus UNIV$ rekursiv und es gibt eine DTM \mathcal{A}_0 , die \overline{UNIV} entscheidet.

Wir betrachten nun die Sprache

$$D = \{\text{code}(\mathcal{A}) \mid \text{code}(\mathcal{A})\text{code}(\mathcal{A}) \notin UNIV\}$$

(d.h. die Maschine \mathcal{A} akzeptiert ihre eigene Kodierung nicht).

Diese Menge kann leicht mit Hilfe von \mathcal{A}_0 entschieden werden:

- Bei Eingabe x dupliziert man x und
- startet dann \mathcal{A}_0 mit Eingabe xx .

Es sei \mathcal{A}_D die DTM, die D entscheidet. Es gilt nun (für $\mathcal{A} = \mathcal{A}_D$):

$$\begin{aligned} \mathcal{A}_D \text{ akzeptiert } \text{code}(\mathcal{A}_D) & \text{ gdw. } \text{code}(\mathcal{A}_D) \in D \\ & \text{gdw. } \text{code}(\mathcal{A}_D)\text{code}(\mathcal{A}_D) \notin UNIV \\ & \text{gdw. } \mathcal{A}_D \text{ akzeptiert } \text{code}(\mathcal{A}_D) \text{ nicht.} \end{aligned}$$

Widerspruch. □

Die in Teil 2 des Beweises verwendete Vorgehensweise nennt man *Diagonalisierung*. Wir hatten hier ja zwei Dimensionen

- Kodierung der TM
- Eingabe der TM,

und die Menge D identifiziert die beiden (betrachtet TM, die als Eingabe ihre eigene Kodierung haben), d.h. D bewegt sich in der Diagonalen.

Ein anderes Beispiel für ein Diagonalisierungsargument ist Cantors Beweis für die Überabzählbarkeit von \mathbb{R} .

Wir wenden im folgenden Einschub Diagonalisierung an, um zu zeigen, dass es nicht-kontextsensitive Typ-0-Sprachen gibt.

Satz 16.5

$$\mathcal{L}_1 \subset \mathcal{L}_0$$

Beweis. Es sei

- G_0, G_1, \dots eine effektive (d.h. mit TM machbare) Aufzählung aller kontextsensitiven Grammatiken mit Terminalalphabet $\Sigma = \{a, b\}$
- und w_0, w_1, \dots eine effektive Aufzählung aller Wörter über Σ .

Wir definieren nun $L \subseteq \{a, b\}^*$ durch

$$w_i \in L \Leftrightarrow w_i \notin L(G_i) \text{ (Diagonalisierung!).}$$

- Es ist L Turing-akzeptierbar, da man für eine kontextsensitive Grammatik G_i und ein Wort w_i entscheiden kann, ob $w_i \in L(G_i)$ ist.

Beachte:

nichtkürzende Produktionen \Rightarrow Aufzählen aller ableitbaren Wörter bis zur Länge $|w_i|$ genügt.

- L ist nicht kontextsensitiv. Anderenfalls gäbe es einen Index k mit $L = L(G_k)$. Nun ist aber

$$w_k \in L(G_k) \text{ gdw. } w_k \in L \text{ gdw. } w_k \notin L(G_k).$$

Widerspruch.

□

Aus der Unentscheidbarkeit von $UNIV$ kann man weitere wichtige Unentscheidbarkeitsresultate herleiten.

Satz 16.6

Das Wortproblem für DTM (und damit auch für \mathcal{L}_0) ist unentscheidbar, d.h. es gibt kein Berechnungsverfahren, das zu jeder gegebenen DTM \mathcal{A} und jedem Eingabewort w entscheidet, ob \mathcal{A} das Wort w akzeptiert.

Beweis. Wenn das Wortproblem für DTM entscheidbar wäre, so wäre auch $UNIV$ entscheidbar (ist es aber nicht).

Um zu entscheiden, ob $\text{code}(\mathcal{A})w \in UNIV$ ist, müsste man ja nur das Entscheidungsverfahren für das Wortproblem feststellen lassen, ob \mathcal{A} das Wort w akzeptiert. □

Satz 16.7

Das Halteproblem für DTM ist unentscheidbar, d.h. es gibt kein Berechnungsverfahren, das zu jeder gegebenen DTM \hat{A} entscheidet, ob \hat{A} beginnend mit leerem Eingabeband terminiert.

Beweis. Wäre das Halteproblem entscheidbar, so auch das Wortproblem.

Um zu gegebener TM \mathcal{A} und gegebenem Wort w zu entscheiden, ob \mathcal{A} das Wort w akzeptiert, könnte man dann nämlich wie folgt vorgehen:

Modifiziere \mathcal{A} zu einer TM $\hat{\mathcal{A}}$:

- Bei leerem Band schreibt $\hat{\mathcal{A}}$ zunächst w auf das Band.
- Danach verhält sich $\hat{\mathcal{A}}$ wie \mathcal{A} .
- Stoppt \mathcal{A} mit *akzeptierender* Stoppkonfiguration, so stoppt auch $\hat{\mathcal{A}}$.
Hält \mathcal{A} mit *nichtakzeptierender* Stoppkonfiguration, so geht $\hat{\mathcal{A}}$ in eine Endlosschleife.

Damit gilt:

$\hat{\mathcal{A}}$ hält mit leerem Eingabeband gdw. \mathcal{A} akzeptiert w .

Mit dem Entscheidungsverfahren für das Halteproblem (angewandt auf $\hat{\mathcal{A}}$) könnte man also das Wortproblem („Ist w in $L(\mathcal{A})$?“) entscheiden. □

Satz 16.8

Das Leerheitsproblem für DTM (und damit auch für \mathcal{L}_0) ist unentscheidbar, d.h. es gibt kein Berechnungsverfahren, das bei gegebener DTM $\hat{\mathcal{A}}$ entscheidet, ob es eine Eingabe w gibt, auf der $\hat{\mathcal{A}}$ terminiert.

Beweis. Wäre das Leerheitsproblem entscheidbar, so auch das Halteproblem.

Um bei gegebener DTM \mathcal{A} zu entscheiden, ob \mathcal{A} auf leerer Eingabe hält, konstruiert man die DTM $\hat{\mathcal{A}}$ wie folgt:

- $\hat{\mathcal{A}}$ löscht seine Eingabe.
- Danach verhält sich $\hat{\mathcal{A}}$ wie \mathcal{A} .
- Stoppt die Berechnung, so geht $\hat{\mathcal{A}}$ in eine akzeptierende Stoppkonfiguration.

Offenbar gibt es eine Eingabe, für die $\hat{\mathcal{A}}$ hält gdw. \mathcal{A} auf leerem Eingabeband hält. □

Satz 16.9

Das Äquivalenzproblem für DTM (und damit auch für \mathcal{L}_0) ist unentscheidbar.

Beweis. Offenbar kann man leicht eine DTM $\hat{\mathcal{A}}$ konstruieren mit $L(\hat{\mathcal{A}}) = \emptyset$.

Wäre das Äquivalenzproblem

$$L(\mathcal{A}_1) = L(\mathcal{A}_2)?$$

entscheidbar, so könnte man durch den Test

$$L(\mathcal{A}) = L(\hat{\mathcal{A}})?$$

das Leerheitsproblem für \mathcal{A} entscheiden. □

Satz 16.10

\mathcal{L}_0 ist nicht unter Komplement abgeschlossen.

Beweis. Wir wissen von der in Satz 16.4 eingeführten Sprache $UNIV$:

- $UNIV$ ist partiell entscheidbar, d.h. gehört zu \mathcal{L}_0 .
- $UNIV$ ist *nicht* entscheidbar.

Wäre $\overline{UNIV} \in \mathcal{L}_0$, d.h. partiell entscheidbar, so würde aber mit Satz 13.4 (Teil 4) folgen, dass $UNIV$ *entscheidbar* ist. □

Das in den Beweisen der Sätze 16.6 bis 16.9 gewählte Vorgehen nennt man *Reduktion*:

- Ein Problem P_1 (z.B. Halteproblem) wird auf ein Problem P_2 (z.B. Äquivalenzproblem) reduziert.
- Wäre daher P_2 entscheidbar, so auch P_1 .
- Weiss man bereits, dass P_1 unentscheidbar ist, so folgt daher, dass auch P_2 unentscheidbar ist.

Formaler betrachten wir (o.B.d.A.) einstellige Relationen $R \subseteq \Sigma^*$.

Definition 16.11 (Reduktion)

- 1) Eine *Reduktion* von $R_1 \subseteq \Sigma^*$ auf $R_2 \subseteq \Sigma^*$ ist eine berechenbare Funktion

$$f : \Sigma^* \rightarrow \Sigma^*,$$

für die gilt:

$$w \in R_1 \quad \text{gdw.} \quad f(w) \in R_2.$$

- 2) Wir schreiben

$$R_1 \leq_m R_2 \quad (R_1 \text{ wird auf } R_2 \text{ reduziert}),$$

falls es eine Reduktion von R_1 nach R_2 gibt.

Lemma 16.12

- 1) $R_1 \leq_m R_2$ und R_2 *entscheidbar* $\Rightarrow R_1$ *entscheidbar*.

2) $R_1 \leq_m R_2$ und R_1 unentscheidbar $\Rightarrow R_2$ unentscheidbar.

Beweis.

1) Um „ $w \in R_1$ “ zu entscheiden,

- berechnet man $f(w)$ und
- entscheidet „ $f(w) \in R_2$ “.

2) Folgt unmittelbar aus 1). □

Mit Hilfe einer Reduktion des Halteproblems kann man zeigen:

Jede *nichttriviale semantische Eigenschaft* von Programmen (DTM) ist *unentscheidbar*.

- *Semantisch* heißt hier: Die Eigenschaft hängt nicht von der syntaktischen Form des Programms, sondern nur von der berechneten Funktion ab.
- *Nichttrivial*: Es gibt berechenbare Funktionen, die sie erfüllen, aber nicht alle berechenbaren Funktionen erfüllen sie.

Satz 16.13 (Satz von Rice)

Es sei E eine Eigenschaft partiell berechenbarer Funktionen $f : \Sigma^* \rightarrow \Sigma^*$, so dass gilt:

$$\emptyset \subset \{f : \Sigma^* \rightarrow \Sigma^* \mid f \text{ erfüllt } E\} \subset \{f : \Sigma^* \rightarrow \Sigma^* \mid f \text{ ist partiell berechenbar}\}$$

Dann ist

$$L(E) := \{\text{code}(\mathcal{A}) \mid \text{die von } \mathcal{A} \text{ berechnete Funktion } f : \Sigma^* \rightarrow \Sigma^* \text{ erfüllt } E\}$$

unentscheidbar.

Beweis. Angenommen, $L(E)$ ist entscheidbar.

Wir zeigen, dass man dann auch ein Entscheidungsverfahren für das Halteproblem erhält (Widerspruch zu Satz 16.7).

Mit f_u bezeichnen wir die überall undefinierte Funktion, welche offenbar partiell berechenbar ist. O.B.d.A. erfülle f_u die Eigenschaft E .

Sonst könnte man statt E die Eigenschaft \bar{E} : „ f erfüllt E nicht“ betrachten. Mit $L(E)$ ist auch $L(\bar{E}) = \bar{L(E)}$ entscheidbar.

Da nicht alle partiell berechenbaren Funktionen E erfüllen, gibt es eine partiell berechenbare Funktion g , welche E *nicht* erfüllt.

Es sei nun \mathcal{A}_g eine DTM für g und \mathcal{A}_u eine für f_u . Wir konstruieren nun eine DTM \mathcal{A}' , welche eine Funktion $f : \Sigma^* \rightarrow \Sigma^*$ berechnet mit

$$w = \text{code}(\mathcal{A}) \text{ für eine DTM } \mathcal{A}, \quad \text{gdw. } f(w) \notin L(E) \quad (*)$$

die auf leerer Eingabe terminiert

Wäre daher $L(E)$ entscheidbar, so auch das Halteproblem.

Konstruktion von \mathcal{A}' :

Die DTM \mathcal{A}' testet bei Eingabe w zunächst, ob w Kodierung einer DTM ist.

- Falls nein, so gibt \mathcal{A}' $\text{code}(\mathcal{A}_u) \in L(E)$ aus.
- Falls ja, so ist $w = \text{code}(\mathcal{A})$ für eine DTM \mathcal{A} . Die DTM \mathcal{A}' gibt dann $\text{code}(\mathcal{A}'')$ aus, wobei \mathcal{A}'' noch geeignet zu definieren ist:

$$\text{code}(\mathcal{A}'') \notin L(E) \text{ gdw. } \mathcal{A} \text{ hält auf leerer Eingabe}$$

Definition von \mathcal{A}'' :

- \mathcal{A}'' ignoriert zunächst die Eingabe x und simuliert das Verhalten von \mathcal{A} auf dem leeren Eingabeband.
- Im Anschluss (d.h. falls \mathcal{A} auf leerem Eingabeband terminiert) verhält sich \mathcal{A}'' wie \mathcal{A}_g bei Eingabe x .

Damit gilt für \mathcal{A}'' offenbar:

- Terminiert \mathcal{A} auf leerer Eingabe *nicht*, so berechnet \mathcal{A}'' die Funktion f_u , d.h. $\text{code}(\mathcal{A}'') \in L(E)$, da f_u E erfüllt.
- Terminiert \mathcal{A} auf leerer Eingabe, so berechnet \mathcal{A}'' die Funktion g , d.h. $\text{code}(\mathcal{A}'') \notin L(E)$, da g E nicht erfüllt.

Insgesamt haben wir also gezeigt, dass die von \mathcal{A}' berechnete Funktion f die Bedingung (\star) erfüllt, das Halteproblem also auf das Problem, $L(E)$ zu entscheiden, reduziert.

Da das Halteproblem unentscheidbar ist, folgt die Unentscheidbarkeit von $L(E)$. □

17. Weitere unentscheidbare Probleme

Die bisher als unentscheidbar nachgewiesenen Probleme betreffen alle nur (semantische) Eigenschaften von Programmen. Aber auch Probleme, die (direkt) nichts mit Programmen zu tun haben, können unentscheidbar sein.

Das folgende von Emil Post definierte Problem ist sehr nützlich, um mittels Reduktion Unentscheidbarkeit zu zeigen.

Definition 17.1 (Postsches Korrespondenzproblem)

Eine Instanz des *Postschen Korrespondenzproblems (PKP)* ist gegeben durch eine endliche Folge

$$P = (x_1, y_1), \dots, (x_k, y_k)$$

von Wortpaaren mit $x_i, y_i \in \Sigma^+$ für ein endliches Alphabet Σ .

Eine *Lösung* des Problems ist eine *Indexfolge* i_1, \dots, i_m mit

- $m > 0$ und
- $i_j \in \{1, \dots, k\}$,

so dass gilt: $x_{i_1} \dots x_{i_m} = y_{i_1} \dots y_{i_m}$.

Beispiel:

- 1) $P_1 = (a, aaa), (abaa, ab), (aab, b)$
 hat z.B. die Folgen 2, 1 und 1, 3 als Lösungen

$$\begin{array}{ll} abaa|a & a|aab \\ ab|aaa & aaa|b \end{array}$$

und damit auch 2, 1, 1, 3 sowie 2, 1, 2, 1, 2, 1, ...

- 2) $P_2 = (ab, aba), (baa, aa), (aba, baa)$
 hat keine Lösung:

$$\begin{array}{l} ab|aba|aba| \dots \\ aba|baa|baa| \dots \end{array}$$

Um die Unentscheidbarkeit des PKP zu zeigen, führen wir zunächst ein Zwischenproblem ein, das *modifizierte PKP (MPKP)*:

Hier muss für die Lösung zusätzlich $i_1 = 1$ gelten, d.h. das Wortpaar, mit dem man beginnen muss, ist festgelegt.

Lemma 17.2

Das MPKP kann auf das PKP reduziert werden.

Beweis. Es sei $P = (x_1, y_1), \dots, (x_k, y_k)$ eine Instanz des MPKP über dem Alphabet Σ . Es seien #, \$ Symbole, die nicht in Σ vorkommen.

Wir definieren die Instanz $f(P)$ des PKP über $\hat{\Sigma} = \Sigma \cup \{\#, \$\}$ wie folgt:

$$f(P) := (x'_0, y'_0), (x'_1, y'_1), \dots, (x'_k, y'_k), (x'_{k+1}, y'_{k+1}),$$

wobei gilt:

- Für $1 \leq i \leq k$ entsteht x'_i aus x_i , indem man *hinter jedem* Symbol ein $\#$ einfügt. Ist z.B. $x_i = abb$, so ist $x'_i = a\#b\#b\#$.
- Für $1 \leq i \leq k$ entsteht y'_i aus y_i , indem man *vor jedem* Symbol ein $\#$ einfügt.
- $x'_0 := \#x'_1$ und $y'_0 := y'_1$
- $x'_{k+1} := \$$ und $y'_{k+1} := \#\$$

Offenbar ist f berechenbar, und man kann leicht zeigen, dass gilt:

„Das MPKP P hat eine Lösung.“ gdw. „Das PKP $f(P)$ hat eine Lösung.“ □

Beispiel:

$P = (a, aaa), (aab, b)$ ist als MPKP lösbar mit Lösung 1, 2.

$$f(P) = \begin{matrix} (x'_0, y'_0) & (x'_1, y'_1) & (x'_2, y'_2) & (x'_3, y'_3) \\ (\#a\#, \#a\#a\#a), & (a\#, \#a\#a\#a), & (a\#a\#b\#, \#b), & (\$, \#\$) \end{matrix}$$

Die Lösung 1, 2 von P liefert die Lösung 0, 2, 3 von $f(P)$:

$$\begin{array}{l} \#a\#|a\#a\#b\#|\$ \\ \#a\#|a\#a|\#b|\#\$ \end{array}$$

Offenbar muss jede Lösung von $f(P)$ mit 0 beginnen.

Wäre daher das PKP entscheidbar, so auch das MPKP. Um die Unentscheidbarkeit des PKP zu zeigen, genügt es also zu zeigen, dass das MPKP unentscheidbar ist.

Lemma 17.3

Das Halteproblem kann auf das MPKP reduziert werden.

Beweis. Gegeben sei eine DTM $\mathcal{A} = (Q, \Sigma, \Gamma, q_0, \Delta, F)$ und ein Eingabewort $w \in \Sigma^*$.

Wir müssen zeigen, wie man \mathcal{A}, w effektiv in eine Instanz $f(\mathcal{A}, w)$ des MPKP überführen kann, so dass gilt:

$$\mathcal{A} \text{ hält auf Eingabe } w \text{ gdw. } f(\mathcal{A}, w) \text{ hat eine Lösung.}$$

Wir verwenden für das MPKP $f(\mathcal{A}, w)$ das Alphabet $\Gamma \cup Q \cup \{\#\}$ mit $\# \notin \Gamma \cup Q$. Das MPKP $f(\mathcal{A}, w)$ besteht aus den folgenden Wortpaaren:

1) Anfangsregel:

$$(\#, \# \flat q_0 w \flat \#)$$

2) Kopierregeln:

$$(a, a) \text{ für alle } a \in \Gamma \cup \{\#\}$$

3) Übergangsregeln:

$$(qa, q'a') \text{ falls } (q, a, a', n, q') \in \Delta$$

$$(qa, a'q') \text{ falls } (q, a, a', r, q') \in \Delta$$

$$(bqa, q'ba') \text{ falls } (q, a, a', l, q') \in \Delta \text{ und } b \in \Gamma$$

$$(\#qa, \#q'\flat a') \text{ falls } (q, a, a', l, q') \in \Delta$$

$$(q\#, q'a'\#) \text{ falls } (q, \flat, a', n, q') \in \Delta$$

$$(q\#, a'q'\#) \text{ falls } (q, \flat, a', r, q') \in \Delta$$

$$(bq\#, q'ba'\#) \text{ falls } (q, \flat, a', l, q') \in \Delta$$

$$(\#q\#, \#q'\flat a'\#) \text{ falls } (q, \flat, a', l, q') \in \Delta$$

4) Löseregeln:

$$(aq, q) \text{ und } (qa, q) \text{ für alle } a \in \Gamma \text{ und } q \in Q \text{ Stoppzustand}$$

(O.B.d.A. hänge in \mathcal{A} das Stoppen nur vom erreichten Zustand ab.)

5) Abschlussregel:

$$(q\#\#, \#) \text{ für alle } q \in Q \text{ mit } q \text{ Stoppzustand}$$

Falls \mathcal{A} bei Eingabe w hält, so gibt es eine Folge von Konfigurationsübergängen

$$k_0 \vdash_{\mathcal{A}} k_1 \vdash_{\mathcal{A}} \dots \vdash_{\mathcal{A}} k_t$$

mit $k_0 = \flat q_0 w \flat$ und $k_t = u\hat{q}v$ mit \hat{q} Stoppzustand.

Daraus kann man nun eine Lösung des MPKP bauen. Zunächst erzeugt man

$$\#k_0\#k_1\#k_2\#\dots\#$$

$$\#k_0\#k_1\#k_2\#\dots\#k_t\#$$

• Dabei beginnt man mit $(\#, \#k_0\#)$.

- Durch Kopierregeln erzeugt man die Teile von k_0 und k_1 , die sich nicht unterscheiden.
- Der Teil, der sich unterscheidet, wird durch die entsprechende Übergangsregel realisiert.

z.B. $(q, a, a', r, q') \in \Delta$ und $k_0 = \#qab\#$

$\#| \#| qa| b| \#| \#|$
 $\# \# qa b \# \#| \#| a'q'| b| \#| \#|$

Man erhält so:

$\#k_0\#$
 $\#k_0\#k_1\#$

- Nun macht man dies so weiter, bis die Stoppkonfiguration k_t mit Stoppzustand \hat{q} erreicht ist. Durch Verwenden von Löschregeln und Kopierregeln löscht man nacheinander die dem Stoppzustand benachbarten Symbole von k_t , z.B.:

$\dots \#a\hat{q}| b| \#| \hat{q} b| \#$
 $\dots \#a \hat{q} b \#| \hat{q}| b| \#| \hat{q}| \#$

- Danach wendet man die Abschlussregel an:

$\dots \#| \hat{q} \# \#$
 $\dots \# \hat{q} \#| \#$

Umgekehrt zeigt man leicht, dass jede Lösung des MPKP einer haltenden Folge von Konfigurationsübergängen entspricht, welche mit k_0 beginnt:

- Man muss mit k_0 beginnen, da wir das MPKP betrachten.
- Durch Kopier- und Übergangsregeln kann man die erzeugten Wörter nicht gleich lang machen.
- Daher muss ein Stoppzustand erreicht werden, damit Lösch- und Abschlussregeln eingesetzt werden können.

□

Da das Halteproblem unentscheidbar ist, folgt die Unentscheidbarkeit des MPKP und damit (wegen Lemma 17.2) die Unentscheidbarkeit des PKP.

Satz 17.4

Das PKP ist unentscheidbar.

Wir verwenden dieses Resultat, um Unentscheidbarkeit von Problemen für kontextfreie und kontextsensitive Sprachen nachzuweisen. Wir zeigen zunächst:

Lemma 17.5

Es ist nicht entscheidbar, ob für kontextfreie Grammatiken G_1, G_2 gilt:

$$L(G_1) \cap L(G_2) \neq \emptyset.$$

Beweis. Wir reduzieren das PKP auf dieses Problem, d.h. wir zeigen:

Zu jeder Instanz P des PKP kann man effektiv kontextfreie Grammatiken $G_P^{(l)}, G_P^{(r)}$ konstruieren, so dass gilt:

$$P \text{ hat Lösung} \quad \underline{\text{gdw.}} \quad L(G_P^{(l)}) \cap L(G_P^{(r)}) \neq \emptyset.$$

Da das PKP unentscheidbar ist, folgt dann die Aussage des Lemmas.

Es sei $P = (x_1, y_1), \dots, (x_k, y_k)$. Wir definieren $G_P^{(l)} = (N_l, \Sigma_l, P_l, S_l)$ mit

- $N_l = \{S_l\}$,
- $\Sigma_l = \Sigma \cup \{1, \dots, k\}$ und
- $P_l = \{S_l \longrightarrow w_i S_l i, S_l \longrightarrow w_i i \mid 1 \leq i \leq k\}$.

$G_P^{(r)}$ wird entsprechend definiert. Es gilt:

$$L(G_P^{(l)}) = \{x_{i_1} \dots x_{i_m} i_m \dots i_1 \mid m \geq 1, i_j \in \{1, \dots, k\}\}$$

$$L(G_P^{(r)}) = \{y_{i_1} \dots y_{i_m} i_m \dots i_1 \mid m \geq 1, i_j \in \{1, \dots, k\}\}$$

Daraus folgt nun unmittelbar:

$$L(G_P^{(l)}) \cap L(G_P^{(r)}) \neq \emptyset$$

gdw. $\exists m \geq 1 \exists i_1, \dots, i_m \in \{1, \dots, k\} : x_{i_1} \dots x_{i_m} i_m \dots i_1 = y_{i_1} \dots y_{i_m} i_m \dots i_1$

gdw. P hat Lösung. □

Da jede kontextfreie Sprache auch kontextsensitiv ist und da die kontextsensitiven Sprachen unter Durchschnitt abgeschlossen sind (Satz 12.6), folgt daraus unmittelbar:

Satz 17.6

Für \mathcal{L}_1 sind das Leerheitsproblem und das Äquivalenzproblem unentscheidbar.

Beachte:

Das Leerheitsproblem ist ein Spezialfall des Äquivalenzproblems, da

$$L(G) = \emptyset \text{ gdw. } L(G) = L(G_\emptyset) \quad (G_\emptyset : \text{kontextsensitive Grammatik für } \emptyset).$$

Für kontextfreie Sprachen wissen wir bereits, dass das Leerheitsproblem entscheidbar ist. Dies ist kein Widerspruch zu Lemma 17.5, da \mathcal{L}_2 nicht unter Durchschnitt abgeschlossen ist.

Satz 17.7

Für \mathcal{L}_2 ist das Äquivalenzproblem unentscheidbar.

Beweis.

1. Man kann sich leicht überlegen, dass die Sprachen $L(G_P^{(l)})$ und $L(G_P^{(r)})$ aus dem Beweis von Lemma 17.5 durch deterministische Kellerautomaten akzeptiert werden können.
2. Die von deterministischen Kellerautomaten akzeptierten kontextfreien Sprachen sind nun aber unter Komplement abgeschlossen.

D.h. es gibt auch einen (effektiv berechenbaren) deterministischen Kellerautomaten und damit eine kontextfreie Grammatik für

$$\overline{L(G_P^{(l)})}$$

(siehe z.B. [Wege93], Satz 8.1.3).

3. Es sei \overline{G} die kontextfreie Grammatik mit $L(\overline{G}) = \overline{L(G_P^{(l)})}$. Nun gilt:

$$\begin{aligned} L(G_P^{(l)}) \cap L(G_P^{(r)}) = \emptyset & \text{ gdw. } L(G_P^{(r)}) \subseteq L(\overline{G}) \\ & \text{ gdw. } L(G_P^{(r)}) \cup L(\overline{G}) = L(\overline{G}) \\ & \text{ gdw. } L(G_{\cup}) = L(\overline{G}), \end{aligned}$$

wobei G_{\cup} die effektiv konstruierbare kontextfreie Grammatik für $L(G_P^{(r)}) \cup L(\overline{G})$ ist.

4. Wäre also das Äquivalenzproblem für kontextfreie Grammatiken entscheidbar, so auch

$$L(G_P^{(l)}) \cap L(G_P^{(r)}) = \emptyset$$

und damit das PKP.

□

Als nächstes verwenden wir das PKP, um zu zeigen, dass die Gültigkeit von Formeln der Prädikatenlogik erster Stufe (PL1) unentscheidbar ist.

Satz 17.8

Das PKP kann auf die Gültigkeit von PL1-Formeln reduziert werden.

Beweis. Es sei $P = (x_1, y_1), \dots, (x_k, y_k)$ ein PKP über dem Alphabet $\Sigma = \{a_1, \dots, a_n\}$.

Wir verwenden

- einstellige Funktionssymbole a_1, \dots, a_n ,
- einstellige Funktionssymbole f_1, \dots, f_k ,
- einstellige Funktionssymbole ℓ, r ,
- das Gleichheitssymbol $=$ (welches als Gleichheit interpretiert werden muss),
- ein Konstantensymbol ε .

Ein Wort

$$w = a_{i_1} \dots a_{i_m} \in \Sigma^*$$

wird dargestellt als Term

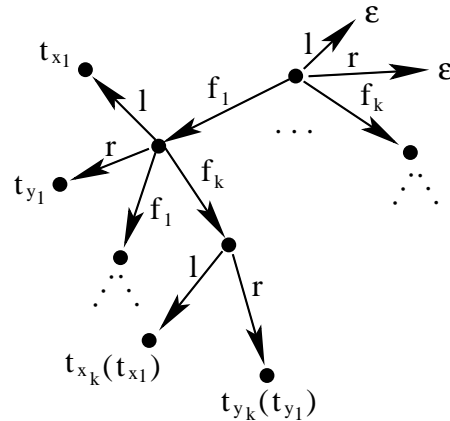
$$t_w = a_{i_m}(a_{i_{m-1}}(\dots(a_{i_1}(\varepsilon))\dots)).$$

Die Konkatenation entspricht daher dem Einsetzen von Termen, d.h.

$$t_{uv} = t_v(t_u).$$

Idee:

Der Suchbaum für eine Lösung des PKP kann wie folgt dargestellt werden:



Dies wird durch die folgende Formel beschrieben:

$$\begin{aligned}
 B_P : & \exists x. (\ell(x) = \varepsilon \quad \wedge \quad r(x) = \varepsilon) \quad \wedge \\
 & \forall x. (\ell(f_1(x)) = t_{x_1}(\ell(x)) \neq \varepsilon \quad \wedge \quad r(f_1(x)) = t_{y_1}(r(x)) \neq \varepsilon) \quad \wedge \\
 & \vdots \\
 & \ell(f_k(x)) = t_{x_k}(\ell(x)) \neq \varepsilon \quad \wedge \quad r(f_k(x)) = t_{y_k}(r(x)) \neq \varepsilon)
 \end{aligned}$$

Die Lösbarkeit des PKP wird dann ausgedrückt durch die Formel:

$$F_P : B_P \Rightarrow \exists z. (\ell(z) = r(z) \neq \varepsilon)$$

Behauptung:

F_P ist gültig gdw. P eine Lösung hat.

denn: Jede Interpretation, die B_P erfüllt, „enthält“ den Suchbaum.

Beweis der Behauptung.

„ \Rightarrow “: Angenommen, F_P ist gültig. Wir benutzen den Suchbaum, um eine Interpretation \mathcal{I} zu definieren, die B_P erfüllt:

- $dom(\mathcal{I}) := \Sigma^* \cup \{1, \dots, k\}^*$
 - $\varepsilon^{\mathcal{I}} := \varepsilon$
 - $a_i : dom(\mathcal{I}) \rightarrow dom(\mathcal{I})$
- $$a_i^{\mathcal{I}}(u) := \begin{cases} ua_i & \text{falls } u \in \Sigma^* \\ \varepsilon & \text{sonst} \end{cases}$$
- $f_j : dom(\mathcal{I}) \rightarrow dom(\mathcal{I})$
- $$f_j^{\mathcal{I}}(u) := \begin{cases} uj & \text{falls } u \in \{1, \dots, k\}^* \\ j & \text{sonst} \end{cases}$$

- $\ell : \text{dom}(\mathcal{I}) \rightarrow \text{dom}(\mathcal{I})$

$$\ell^{\mathcal{I}}(u) := \begin{cases} x_{i_1} \dots x_{i_m} & \text{falls } u = i_1 \dots i_m \in \{1, \dots, k\}^* \\ \varepsilon & \text{sonst} \end{cases}$$
- $r : \text{dom}(\mathcal{I}) \rightarrow \text{dom}(\mathcal{I})$

$$r^{\mathcal{I}}(u) := \begin{cases} y_{i_1} \dots y_{i_m} & \text{falls } u = i_1 \dots i_m \in \{1, \dots, k\}^* \\ \varepsilon & \text{sonst} \end{cases}$$

Man sieht leicht, dass \mathcal{I} die Formel B_P wahr macht:

- $\exists x. \ell(x) = \varepsilon \wedge r(x) = \varepsilon$ ist erfüllt, da

$$\ell^{\mathcal{I}}(\varepsilon^{\mathcal{I}}) = \varepsilon = \varepsilon^{\mathcal{I}} \text{ und } r^{\mathcal{I}}(\varepsilon^{\mathcal{I}}) = \varepsilon = \varepsilon^{\mathcal{I}}.$$

- Sei $u \in \text{dom}(\mathcal{I})$:

$$- u = i_1 \dots i_m \in \{1, \dots, k\}^*:$$

$$\begin{aligned} \ell^{\mathcal{I}}(f_j^{\mathcal{I}}(u)) &= \ell^{\mathcal{I}}(i_1 \dots i_m j) = x_{i_1} \dots x_{i_m} x_j \\ &\quad \parallel \\ t_{x_j}^{\mathcal{I}}(\ell^{\mathcal{I}}(u)) &= t_{x_j}^{\mathcal{I}}(x_{i_1} \dots x_{i_m}) = x_{i_1} \dots x_{i_m} x_j \end{aligned}$$

$$- u \notin \{1, \dots, k\}^*:$$

$$\begin{aligned} \ell^{\mathcal{I}}(f_j^{\mathcal{I}}(u)) &= \ell^{\mathcal{I}}(j) = x_j \\ &\quad \parallel \\ t_{x_j}^{\mathcal{I}}(\ell^{\mathcal{I}}(u)) &= t_{x_j}^{\mathcal{I}}(\varepsilon) = x_j \end{aligned}$$

- Für r geht das entsprechend.

Also erfüllt \mathcal{I} die Formel B_P . Damit muss \mathcal{I} auch

$$\exists z. (\ell(z) = r(z) \neq \varepsilon)$$

erfüllen, d.h. es gibt ein $u \in \text{dom}(\mathcal{I})$ mit

$$\ell^{\mathcal{I}}(u) = r^{\mathcal{I}}(u) \neq \varepsilon^{\mathcal{I}} = \varepsilon.$$

1. Fall: $u = i_1 \dots i_m \in \{1, \dots, k\}^+$. Dann ist

$$\begin{aligned} \ell^{\mathcal{I}}(u) &= x_{i_1} \dots x_{i_m} \\ &\quad \parallel \\ r^{\mathcal{I}}(u) &= y_{i_1} \dots y_{i_m} \end{aligned}$$

d.h. $i_1 \dots i_m$ ist eine Lösung des PKP.

2. Fall: $u \in \Sigma^*$.

- $u = \varepsilon : \ell^{\mathcal{I}}(u) = \varepsilon = r^{\mathcal{I}}(u)$ kann nicht sein
- $u \in \Sigma^+ : \text{Dann ist ebenfalls } \ell^{\mathcal{I}}(u) = \varepsilon = r^{\mathcal{I}}(u), \text{ d.h. auch dieser Fall kann nicht eintreten.}$

„ \Leftarrow “: Angenommen, P hat eine Lösung. Dann sei \mathcal{I} eine Interpretation, die B_P wahr macht.

Zu zeigen: \mathcal{I} erfüllt auch

$$\exists z. (\ell(z) = r(z) \neq \varepsilon).$$

Es sei $i_1 \dots i_m$ für $m > 0$ eine Lösung von P , d.h.

$$x_{i_1} \dots x_{i_m} = y_{i_1} \dots y_{i_m}.$$

Da \mathcal{I} die Formel B_P erfüllt, gilt:

$$\begin{aligned} \ell^{\mathcal{I}}(f_{i_m}^{\mathcal{I}}(\dots f_{i_1}^{\mathcal{I}}(\varepsilon^{\mathcal{I}}))) &= t_{x_{i_1} \dots x_{i_m}}^{\mathcal{I}} \neq \varepsilon^{\mathcal{I}} \\ r^{\mathcal{I}}(f_{i_m}^{\mathcal{I}}(\dots f_{i_1}^{\mathcal{I}}(\varepsilon^{\mathcal{I}}))) &= t_{y_{i_1} \dots y_{i_m}}^{\mathcal{I}} \neq \varepsilon^{\mathcal{I}} \end{aligned}$$

Aus $x_{i_1} \dots x_{i_m} = y_{i_1} \dots y_{i_m}$ folgt $t_{x_{i_1} \dots x_{i_m}}^{\mathcal{I}} = t_{y_{i_1} \dots y_{i_m}}^{\mathcal{I}}$,

d.h. $f_{i_m}^{\mathcal{I}}(\dots f_{i_1}^{\mathcal{I}}(\varepsilon^{\mathcal{I}}))$ ist das gesuchte z .

□

IV. Komplexität

Einführung

In der Praxis genügt es nicht zu wissen, dass eine Funktion berechenbar ist. Man interessiert sich auch dafür, wie groß der *Aufwand* zur Berechnung ist.

Dabei betrachtet man verschiedene Fragestellungen in Bezug auf *Aufwand/Komplexität*:

Rechenzeit

(bei TM: Anzahl der Übergänge bis zum Halten)

Speicherplatz

(bei TM: Anzahl der benutzten Felder)

Beides soll abgeschätzt werden als *Funktion in der Größe der Eingabe*.

Man kann die Komplexität eines speziellen Algorithmus/Programms betrachten (wie in der Vorlesung „Algorithmen und Datenstrukturen“) oder die

Komplexität eines Entscheidungsproblems: Wieviel Aufwand benötigt der „beste“ Algorithmus im „schlimmsten“ Fall?

Wir betrachten hier den zweiten Fall.

18. Komplexitätsklassen

Wir führen im folgenden den Begriff der *Komplexitätsklasse* allgemein ein und definieren dann einige fundamentale *Zeit-* und *Platzkomplexitätsklassen*. Zunächst führen wir formal ein, was es heißt, dass der Aufwand durch eine Funktion der Größe der Eingabe *beschränkt* ist.

Definition 18.1 (*f(n)*-zeitbeschränkt, *f(n)*-platzbeschränkt)

Es sei $f : \mathbb{N} \rightarrow \mathbb{N}$ eine Funktion und \mathcal{A} eine DTM über Σ .

- 1) \mathcal{A} heißt *f(n)*-zeitbeschränkt, falls für alle $w \in \Sigma^*$ die Maschine \mathcal{A} bei Eingabe w nach $\leq f(|w|)$ Schritten anhält.
- 2) \mathcal{A} heißt *f(n)*-platzbeschränkt, falls für alle $w \in \Sigma^*$ die Maschine \mathcal{A} bei Eingabe w $\leq f(|w|)$ viele Felder des Bandes benutzt.

Auf *NTM* kann man die Definition dadurch übertragen, dass die Aufwandsbeschränkung für *alle* bei der gegebenen Eingabe *möglichen Berechnungen* zutreffen muss.

Wir betrachten im folgenden nicht (wie in „Algorithmen und Datenstrukturen“) die Komplexität einzelner Algorithmen (DTMs), sondern die *Komplexität von Entscheidungsproblemen*, wobei wir uns o.B.d.A. auf einstellige Relationen (d.h. $L \subseteq \Sigma^*$) beschränken.

Definition 18.2 (Komplexitätsklassen)

$$\begin{aligned} DTIME(f(n)) &:= \{L \mid \text{es gibt eine } f(n)\text{-zeitbeschränkte DTM, die } L \text{ entscheidet}\} \\ NTIME(f(n)) &:= \{L \mid \text{es gibt eine } f(n)\text{-zeitbeschränkte NTM, die } L \text{ akzeptiert}\} \\ DSPACE(f(n)) &:= \{L \mid \text{es gibt eine } f(n)\text{-platzbeschränkte DTM, die } L \text{ entscheidet}\} \\ NSPACE(f(n)) &:= \{L \mid \text{es gibt eine } f(n)\text{-platzbeschränkte NTM, die } L \text{ akzeptiert}\} \end{aligned}$$

Man kann einige elementare Zusammenhänge zwischen den Komplexitätsklassen feststellen:

Satz 18.3

- 1) $DTIME(f(n)) \subseteq DSPACE(f(n)) \subseteq NSPACE(f(n))$
- 2) $DTIME(f(n)) \subseteq NTIME(f(n))$

Beweis. Offenbar kann man in k Schritten höchstens k Felder benutzen. Außerdem ist eine DTM ja auch eine NTM. □

Wir betrachten die folgenden *fundamentalen* Komplexitätsklassen:

Definition 18.4 ($P, NP, PSPACE$)

$$\begin{aligned}
 P &:= \bigcup_{p \text{ Polynom in } n} DTIME(p(n)) \\
 NP &:= \bigcup_{p \text{ Polynom in } n} NTIME(p(n)) \\
 PSPACE &:= \bigcup_{p \text{ Polynom in } n} NSPACE(p(n))
 \end{aligned}$$

Bei der Definition von $PSPACE$ hätte man genauso gut $DSPACE$ statt $NSPACE$ verwenden können, d.h.

$$PSPACE = \bigcup_{p \text{ Polynom in } n} DSPACE(p(n)),$$

denn:

Wir hatten gesehen, dass man jede NTM durch eine DTM simulieren kann. Nach einem Resultat von **Savitch** kann diese Simulation so gestaltet werden, dass der Platzbedarf der simulierenden DTM quadratisch im Platzbedarf der NTM ist (ohne Beweis).

Bemerkung 18.5.

Die Bedeutung der in Definition 18.4 eingeführten Komplexitätsklassen ergibt sich u.a. aus den folgenden Überlegungen:

- 1) Alle üblichen berechnungsuniversellen Maschinenmodelle lassen sich auf TM in Polynomialzeit simulieren, d.h.:

Benötigt das Maschinenmodell zur Berechnung n Schritte, so benötigt die simulierende TM $q(n)$ Schritte, wobei q ein Polynom ist.

Es folgt, dass die in Definition 18.4 eingeführten Komplexitätsklassen unabhängig vom speziellen Maschinenmodell sind.

- 2) Probleme, die in P sind, werden i.a. als *effizient lösbar* („tractable“, d.h. „machbar“) bezeichnet. Dies liegt daran, dass Polynome im Vergleich zu Exponentialfunktionen ($n \mapsto 2^n$) ein recht moderates Wachstum haben.

Aber:

In Anwendungen, bei denen man für sehr große Eingaben sehr schnell (in Echtzeit) Antworten haben will, ist ein Polynom mit hohem Grad häufig auch nicht mehr tolerierbar.

- 3) NP enthält viele Probleme, für die derzeit nur deterministische Entscheidungsverfahren mit exponentiellem Zeitaufwand bekannt sind. Man bezeichnet diese Probleme daher häufig als *nicht effizient lösbar* („intractable“).

Aber:

Die eingeführten Komplexitätsklassen (wie NP) betrachten stets den schwierigsten Fall („*worst-case*“-Komplexität). Es kann durchaus sein, dass es wegen gewisser pathologischer Situationen keine polynomiale Schranke für die Rechenzeit gibt, aber für typische Fälle oder im Mittel doch polynomiales Verhalten erzielt werden kann.

Wie ist das Verhältnis zwischen den Komplexitätsklassen? Aus Satz 18.3 ergibt sich:

$$P \subseteq NP \subseteq PSPACE$$

Ob diese Inklusionen *echt* sind, ist bisher nicht bekannt. Die Frage, ob P gleich NP ist, d.h. das

$$P = NP \text{ - Problem,}$$

ist eines der fundamentalen offenen Probleme der Informatik, dessen Beantwortung weitreichende Konsequenzen hätte.

19. NP-vollständige Probleme

Es geht hier um Probleme in NP , die „mindestens so schwer“ wie jedes Problem in NP sind. Ein solches Problem ist SAT:

Definition 19.1 (SAT)

SAT ist die Menge der erfüllbaren („satisfiable“) Booleschen Ausdrücke

- mit Variablen x_i ($i \geq 1$, kodiert als Dualzahl bei Eingabe für TM)
- und Operatoren \neg, \wedge, \vee .

Beispiel:

- Der Ausdruck $(x_1 \wedge x_2) \vee \neg x_1$ ist erfüllbar, da er von der Belegung $x_1 \mapsto 1, x_2 \mapsto 1$ wahr gemacht wird.
- $x_1 \wedge \neg x_1$ ist *nicht* erfüllbar.

Lemma 19.2

$SAT \in NP$

Beweis. Eine NTM akzeptiert die Elemente von SAT wie folgt:

- 1) Teste, ob die Eingabe ein Boolescher Ausdruck ist (entspricht dem Wortproblem für kontextfreie Grammatiken, polynomial entscheidbar).
- 2) Durchlaufe den Ausdruck und ersetze *nichtdeterministisch* jede Variable x_i durch 0 oder 1 (polynomialer Aufwand).
- 3) Werte den Ausdruck aus und akzeptiere, falls das Ergebnis 1 ist (geht polynomial). □

Inwiefern ist nun SAT „das schwierigste“ NP-Problem?

Definition 19.3 (polynomialzeitberechenbar,-reduzierbar, NP-vollständig)

- 1) Die Funktion $f : \Sigma^* \rightarrow \Sigma^*$ heißt *polynomialzeitberechenbar*, falls es
 - ein Polynom p und
 - eine $p(n)$ -zeitbeschränkte DTM gibt,
 - die f berechnet.
- 2) $L \subseteq \Sigma^*$ ist *polynomialzeitreduzierbar* auf $L' \subseteq \Sigma^*$ (geschrieben als $L \leq_p L'$), falls es eine *polynomialzeitberechenbare* Funktion f gibt mit

$$w \in L \quad \text{gdw.} \quad f(w) \in L'$$

für alle $w \in \Sigma^*$.

3) L_0 heißt *NP-vollständig*, falls gilt:

- $L_0 \in NP$ (L_0 ist in NP)
- Für alle $L \in NP$ gilt: $L \leq_p L_0$
 L_0 ist *NP-hart*, d.h. mindestens so hart zu lösen wie jedes andere NP-Problem.

Satz 19.4

1) Ist L_0 *NP-vollständig*, so gilt:

$$L_0 \in P \Rightarrow P = NP$$

(Es ist bisher kein *NP-vollständiges* L_0 bekannt, das in P ist!)

2) *SAT* ist *NP-vollständig*.

Beweis.

1) Es sei L_0 *NP-vollständig* und $L_0 \in P$, d.h. L_0 wird von einer $p(n)$ -zeitbeschränkten DTM (für ein Polynom p) entschieden.

Wir müssen zeigen, dass daraus $NP \subseteq P$ folgt.

Sei also $L \in NP$. Es ist daher $L \leq_p L_0$, d.h. es gibt ein f , das mit Zeitaufwand $\leq q(n)$ berechenbar ist (für Polynom q), so dass

$$w \in L \text{ gdw. } f(w) \in L_0.$$

Die DTM, welche L entscheidet, geht wie folgt vor:

- Bei Eingabe w berechnet sie $f(w)$. Sie benötigt $\leq q(|w|)$ viel Zeit. Daher ist auch die erzeugte Ausgabe $\leq |w| + q(|w|)$.
- Wende Entscheidungsverfahren für L_0 auf $f(w)$ an.

Insgesamt benötigt man somit

$$\leq q(|w|) + p(|w| + q(|w|))$$

viele Schritte, was ein Polynom in $|w|$ ist.

2) Wir beschreiben nur kurz die Idee:

Zu zeigen ist:

$$\forall L : L \in NP \Rightarrow L \leq_p SAT.$$

Sei \mathcal{A} eine polynomialzeitbeschränkte NTM, welche $L \subseteq \Sigma^*$ akzeptiert. Die gesuchte Reduktionsfunktion f weist jedem Wort $w \in \Sigma^*$ einen Booleschen Ausdruck β_w zu, so dass

- Der Übergang von w zu β_w ist in Polynomialzeit berechenbar.
- \mathcal{A} akzeptiert w gdw. β_w ist erfüllbar.

Die genaue Definition von β_w ist sehr aufwendig. Es muss ja die Arbeitsweise von \mathcal{A} auf Eingabe w durch einen Booleschen Ausdruck beschrieben werden.

Idee:

- Aussagenvariablen $Band_{t,i,a}$ für
„Zum Zeitpunkt t steht auf Feld i das Symbol a “.
- Da \mathcal{A} polynomialzeitbeschränkt ist, muss man nur polynomial viele Bandzellen und Zeitpunkte betrachten!
- Aussagenvariablen $Kopf_{t,i}$
„Kopf zum Zeitpunkt t auf Feld i “,
- $Zustand_{t,q}$
„Zustand zum Zeitpunkt t ist q “.
- Übergänge: (p, a, a', r, q) liefert z.B.:
 $Band_{t,i,a} \wedge Kopf_{t,i} \wedge Zustand_{t,p} \longrightarrow Band_{t+1,i,a'} \wedge Kopf_{t+1,i+r} \wedge Zustand_{t+1,q}$.
- Kodierung der Eingabe w zum Zeitpunkt 0.
- Kodierungsbedingungen, z.B.

$$\neg(Band_{t,i,a} \wedge Band_{t,i,b}) \text{ für } a \neq b$$

- etc. □

Weitere NP-vollständige Probleme erhält man durch polynomielle Reduktion.

Satz 19.5

- 1) Ist $L_2 \in NP$ und gilt $L_1 \leq_p L_2$, so ist auch L_1 in NP .
- 2) Ist L_1 NP-hart und gilt $L_1 \leq_p L_2$, so ist auch L_2 NP-hart.

Beweis.

- 1) Wegen $L_2 \in NP$ gibt es eine polynomialzeitbeschränkte NTM \mathcal{A} , welche L_2 akzeptiert.
Wegen $L_1 \leq_p L_2$ gibt es eine polynomialzeitberechenbare Funktion f mit

$$w \in L_1 \text{ gdw. } f(w) \in L_2.$$

Die polynomialzeitbeschränkte NTM für L_2 arbeitet wie folgt:

- Bei Eingabe w berechnet sie zunächst $f(w)$.
- Dann wendet sie \mathcal{A} auf $f(w)$ an.

2) Wir müssen zeigen, dass für alle $L \in NP$ gilt:

$$L \leq_p L_2.$$

Die polynomialzeitberechenbare Reduktionsfunktion f mit

$$w \in L \text{ gdw. } f(w) \in L_2$$

erhält man wie folgt:

- Da L_1 NP-hart ist, gibt es eine polynomialzeitberechenbare Funktion g mit

$$w \in L \text{ gdw. } g(w) \in L_1.$$

- Wegen $L_1 \leq_p L_2$ gibt es eine polynomialzeitberechenbare Funktion h mit

$$u \in L_1 \text{ gdw. } h(u) \in L_2.$$

Wir definieren

$$f(w) := h(g(w)).$$

Dann gilt:

$$w \in L \text{ gdw. } g(w) \in L_1 \text{ gdw. } h(g(w)) \in L_2. \quad \square$$

Zunächst zeigen wir, dass auch ein Teilproblem von *SAT*, bei dem man nur Boolesche Ausdrücke einer ganz speziellen Form zulässt, bereits NP-hart ist.

Definition 19.6 (3SAT)

- Eine *3-Klausel* ist von der Form

$$l_1 \vee l_2 \vee l_3,$$

wobei l_i eine Variable oder eine negierte Variable ist.

- Ein *3SAT-Problem* ist eine endliche Konjunktion von 3-Klauseln.
- *3SAT* ist die Menge der erfüllbaren *3SAT*-Probleme.

Satz 19.7

3SAT ist NP-vollständig.

Beweis.

- 1) $3SAT \in NP$ folgt unmittelbar aus $SAT \in NP$, da jedes *3SAT*-Problem ein Boolescher Ausdruck ist.

- 2) Es sei F ein beliebiger Boolescher Ausdruck. Wir geben ein polynomielles Verfahren an, das F in ein $3SAT$ -Problem F' umwandelt, so dass gilt:

$$F \text{ erfüllbar} \quad \text{gdw.} \quad F' \text{ erfüllbar.}$$

Beachte:

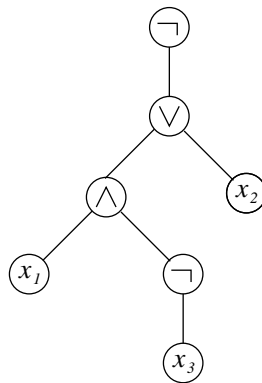
Es ist nicht nötig, dass F und F' (als Boolesche Ausdrücke) äquivalent sind.

Die Umformung erfolgt in mehreren Schritten, die wir am Beispiel des Ausdrucks

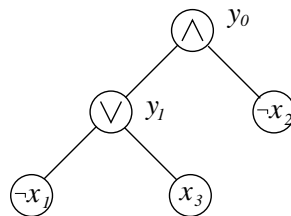
$$\neg((x_1 \wedge \neg x_3) \vee x_2)$$

veranschaulichen.

Wir stellen diesen Ausdruck als Baum dar:



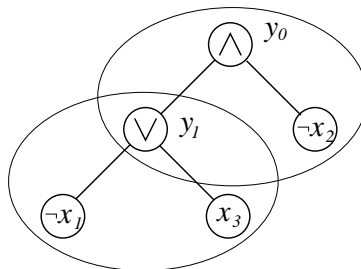
- 1. Schritt:** Wende *de Morgan* an, um die Negationszeichen zu den Blättern zu schieben.



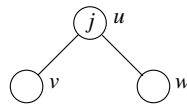
Dies erfordert nur einen Durchlauf durch den Baum, ist also in *linearer* Zeit realisierbar.

- 2. Schritt:** Ordne jedem inneren Knoten eine neue Variable aus $\{y_0, y_1, \dots\}$ zu, wobei die Wurzel y_0 erhält.

- 3. Schritt:** Fasse jede Verzweigung (gedanklich) zu einer Dreiergruppe zusammen:



Jeder Verzweigung der Form



mit $j \in \{\wedge, \vee\}$ ordnen wir eine Formel der folgenden Form zu:

$$(u \Leftrightarrow (v \ j \ w))$$

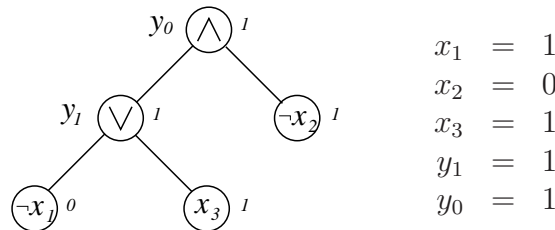
Diese Formeln werden nun konjunktiv mit y_0 verknüpft, was die Formel F_1 liefert.

Im Beispiel ist F_1 :

$$y_0 \wedge (y_0 \Leftrightarrow (y_1 \wedge \neg x_2)) \wedge (y_1 \Leftrightarrow (\neg x_1 \vee x_3))$$

Die Ausdrücke F und F_1 sind *erfüllbarkeitsäquivalent*, denn:

Eine erfüllende Belegung für F kann zu einer für F_1 erweitert werden, indem man die Werte für die Variablen y_i durch die Auswertung der entsprechenden Teilformel bestimmt, z.B.:



Umgekehrt ist jede erfüllende Belegung für F_1 auch eine für F .

4. Schritt: Jede Konjunktion in F_1 wird separat in eine konjunktive Normalform umgeformt:

$$\begin{aligned} a \Leftrightarrow (b \vee c) &\rightsquigarrow (\neg a \vee (b \vee c)) \wedge (\neg(b \vee c) \vee a) \\ &\rightsquigarrow (\neg a \vee b \vee c) \wedge (\neg b \vee a) \wedge (\neg c \vee a) \\ &\rightsquigarrow (\neg a \vee b \vee c) \wedge (\neg b \vee a \vee a) \wedge (\neg c \vee a \vee a) \end{aligned}$$

$$\begin{aligned} a \Leftrightarrow (b \wedge c) &\rightsquigarrow (\neg a \vee (b \wedge c)) \wedge (\neg(b \wedge c) \vee a) \\ &\rightsquigarrow (\neg a \vee b) \wedge (\neg a \vee c) \wedge (\neg b \vee \neg c \vee a) \\ &\rightsquigarrow (\neg a \vee b \vee b) \wedge (\neg a \vee c \vee c) \wedge (\neg b \vee \neg c \vee a) \end{aligned}$$

Dadurch erhält man Formeln der gewünschten *3SAT*-Form.

Beachte:

Pro Teilformel ist der Aufwand hierfür konstant. □

Ein weiteres interessantes NP-vollständiges Problem ist das *Mengenüberdeckungsproblem*.

Definition 19.8 (Mengenüberdeckung)

Gegeben: Ein *Mengensystem* über einer endlichen Grundmenge M , d.h.

$$T_1, \dots, T_k \subseteq M$$

sowie eine Zahl $n \leq k$.

Frage: Gibt es eine Auswahl von n Mengen, die ganz M überdecken, d.h.

$$\{i_1, \dots, i_n\} \subseteq \{1, \dots, k\} \text{ mit } T_{i_1} \cup \dots \cup T_{i_n} = M.$$

Satz 19.9

Mengenüberdeckung ist NP-vollständig.

Beweis.

1) Mengenüberdeckung ist in *NP*:

- Wähle nichtdeterministisch Indices i_1, \dots, i_n und
- überprüfe, ob $T_{i_1} \cup \dots \cup T_{i_n} = M$ gilt.

2) Um NP-Härte zu zeigen, reduzieren wir *3SAT* auf Mengenüberdeckung.

Sei also $F = K_1 \wedge \dots \wedge K_m$ eine Instanz von *3SAT*, welche die Variablen x_1, \dots, x_n enthält.

Wir definieren $M := \{1, \dots, m, m+1, \dots, m+n\}$.

Für jedes $i \in \{1, \dots, n\}$ sei

$$T_i := \{j \mid x_i \text{ kommt in } K_j \text{ als Disjunkt vor}\} \cup \{m+i\}$$

$$T'_i := \{j \mid \neg x_i \text{ kommt in } K_j \text{ als Disjunkt vor}\} \cup \{m+i\}$$

Wir betrachten das Mengensystem:

$$T_1, \dots, T_n, T'_1, \dots, T'_n,$$

und fragen, ob es eine Überdeckung von M mit n Mengen gibt.

a) Für eine gegebene Belegung der Variablen, welche F erfüllt, wählen wir

- T_i , falls $x_i = 1$
- T'_i , falls $x_i = 0$

Dies liefert n Mengen, in denen jedes Element von M vorkommt:

- j mit $1 \leq j \leq m$, da jedes K_j zu 1 evaluiert wird.
- $m+i$ mit $1 \leq i \leq n$, da für jedes i entweder T_i oder T'_i gewählt wird.

b) Sei umgekehrt

$$\{U_1, \dots, U_n\} \subseteq \{T_1, \dots, T_n, T'_1, \dots, T'_n\} \text{ mit } U_1 \cup \dots \cup U_n = M$$

gegeben.

Da für $1 \leq i \leq n$ das Element $m+i$ in $U_1 \cup \dots \cup U_n$ vorkommt, kommt T_i oder T'_i in $\{U_1, \dots, U_n\}$ vor. Da wir nur n verschiedene Mengen haben, kommt jeweils genau eines von beiden vor.

O.B.d.A. sei $U_i \in \{T_i, T'_i\}$ für $i = 1, \dots, n$.

Wir definieren nun die Belegung:

$$x_i = \begin{cases} 1 & \text{falls } U_i = T_i \\ 0 & \text{falls } U_i \neq T_i \end{cases}$$

Diese erfüllt jedes K_j , da j in $U_1 \cup \dots \cup U_n$ vorkommt. □

Viele Probleme für Graphen sind NP-vollständig, z.B. auch das folgende:

Definition 19.10 (Clique)

Gegeben: Ein gerichteter Graph $G = (V, E)$ und eine Zahl $k \in \mathbb{N}$

Frage: Besitzt G eine k -Clique, d.h. eine Teilmenge $U \subseteq V$ mit

- $|U| = k$ und
- für alle $u \neq v$ in U gilt: $(u, v) \in E$.

Satz 19.11

Clique ist NP-vollständig.

Beweis.

1) Clique ist in NP:

- Rate eine Teilmenge der Größe k und
- teste, ob sie eine k -Clique ist.

2) Clique ist NP-hart:

Wir reduzieren 3SAT auf Clique.

Sei also

$$F = \overbrace{(l_{11} \vee l_{12} \vee l_{13})}^{K_1} \wedge \dots \wedge \overbrace{(l_{m1} \vee l_{m2} \vee l_{m3})}^{K_m}$$

mit $l_{ij} \in \{x_1, \dots, x_n\} \cup \{\neg x_1, \dots, \neg x_n\}$.

Der Graph G wird nun wie folgt definiert:

- $V := \{(i, j) \mid 1 \leq i \leq m \text{ und } 1 \leq j \leq 3\}$

- $E := \{((i, j), (p, q)) \mid i \neq p \text{ und } l_{ij} \neq \bar{l}_{pq}\}$, wobei

$$\bar{l} = \begin{cases} \neg x & \text{falls } l = x \\ x & \text{falls } l = \neg x \end{cases} .$$

- $k = m$

Es gilt nun:

F ist erfüllbar mit einer Belegung B

gdw. es gibt in jeder Klausel K_i ein Literal l_{ij_i} mit Wert 1

gdw. es gibt Literale $l_{1j_1}, \dots, l_{mj_m}$, die paarweise nicht komplementär sind

gdw. es gibt Knoten $(1, j_1), \dots, (m, j_m)$, die paarweise miteinander verbunden sind

gdw. es gibt eine k -Clique in G

□

Bemerkung:

Entsprechend zur Definition von NP-Vollständigkeit eines Problems kann man auch PSPACE-Vollständigkeit definieren. Das Äquivalenzproblem für NEAs/reguläre Ausdrücke ist PSPACE-vollständig.

Abkürzungsverzeichnis

bzw.	beziehungsweise
DEA	deterministischer endlicher Automat
d.h.	das heißt
DTM	deterministische Turingmaschine
EBNF	erweiterte Backus-Naur-Form
etc.	et cetera
gdw.	genau dann wenn
geg.	gegeben
i.a.	im allgemeinen
LBA	linear beschränkter Automat
MPKP	modifiziertes Postsches Korrespondenzproblem
NEA	nichtdeterministischer endlicher Automat
NP	nichtdeterministisch polynomiell
NTM	nichtdeterministische Turingmaschine
o.B.d.A.	ohne Beschränkung der Allgemeingültigkeit
PDA	pushdown automaton (Kellerautomat)
PKP	Postsches Korrespondenzproblem
PL1	Prädikatenlogik erster Stufe
SAT	satisfiability problem (Erfüllbarkeitstest der Aussagenlogik)
TM	Turingmaschine (allgemein)
u.a.	unter anderem
URM	unbeschränkte Registermaschine
vgl.	vergleiche
z.B.	zum Beispiel
□	was zu beweisen war (q.e.d)

Literaturverzeichnis

- [Koz06] Dexter Kozen. *Automata and Computability*. Springer Verlag, 2007
- [Hop06] John E. Hopcroft, Rajeev Motwani, Jeffrey D. Ullmann. *Introduction to Automata Theory, Languages, and Computation*. Dritte Ausgabe. Addison Wesley, 2006
- [Schö97] Uwe Schöning. *Theoretische Informatik – kurzgefaßt*. Spektrum Akademischer Verlag, 1997
- [Wege93] Ingo Wegener. *Theoretische Informatik*. Teubner-Verlag, 1993