

Einführung in die Computerlinguistik

Automaten

Hinrich Schütze

Center for Information and Language Processing

2018-12-10

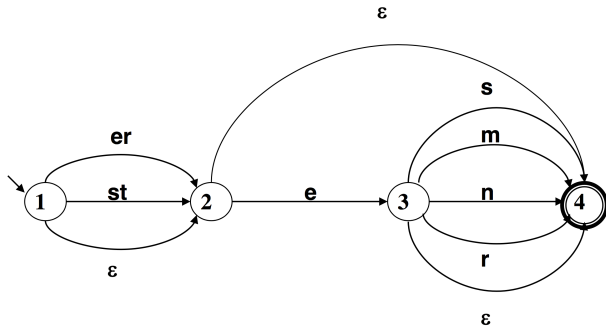
Die Grundfassung dieses Foliensatzes wurde von Prof. Dr. Stefan Evert erstellt. Fehler und Mängel sind ausschließlich meine Verantwortung.

- 1 Motivation
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Automaten
- 5 Breadth-First & Depth-First
- 6 Moodle

- 1 Motivation
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Automaten
- 5 Breadth-First & Depth-First
- 6 Moodle

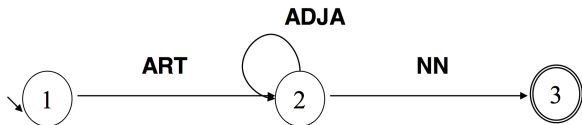
- Formen von “klein” :
klein, kleine, kleinem, kleineren, kleiner, kleines,
kleiner, kleinere, kleinerem, kleinereren, kleinerer, kleineres
kleinst, kleinste, kleinstem, kleinsteren, kleinsten, kleinstes
- Erkennung: Ist eine gegebene Form eine Form von “klein” ?
- Generierung: Generiere für eine gegebene Spezifikation (z.B., stark, Nominativ, Singular, Maskulinum) eine Form (“kleiner”)
- Beide Aufgaben können mit Automaten gelöst werden.

Automat für Adjektive (Buchstabenebene)



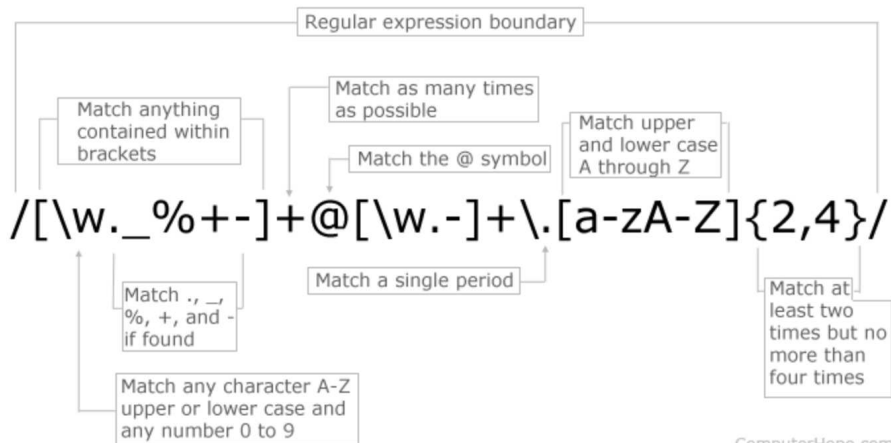
- Ein wichtiger Teil der Syntax natürlicher Sprache kann durch Automaten dargestellt werden.
- Beispiel:
 - das Auto
 - das rote Auto
 - das alte rote Auto
 - das schnelle alte rote Auto
 - das kaputte schnelle alte rote Auto

Automat für Nominalphrase (Lexemformebene)



- Die meisten Textverarbeitungssysteme verwenden Reguläre Ausdrücke.
- Vielfältige Verwendungsmöglichkeiten
- Entfernen von Markup (HTML etc)
- Extraktion von Information
- Automatischen Ersetzen
- ...
- Reguläre Ausdrücke ähneln einer eigenen Programmiersprache (sind aber auch Teil fast jeder Programmiersprache).

Regulärer Ausdruck für Email-Adressen



ComputerHope.com

Theoretische Motivation: Was ist Sprache?

- Ist Sprache endlich?
- Ist das Vokabular einer Sprache begrenzt?
- Können Sätze beliebig lang sein?
- Ist Sprache rekursiv?
- Ist Sprache hierarchisch?
- Die Theorie der Formalen Sprachen hilft bei der Präzisierung und Beantwortung dieser Fragen.
- Zunächst nur Form, nicht Funktion/Bedeutung.
- Form: Kann man das so sagen?
- Funktion: Was ist die Bedeutung des Gesagten?

- Morphologie
- Syntax
- Praktische Motivation: Reguläre Ausdrücke
- Theoretische Motivation: Was ist Sprache?

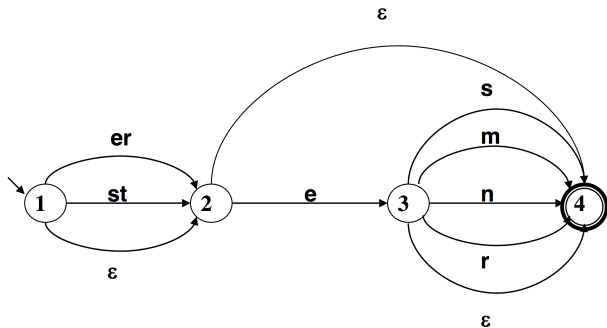
Umfrage

Für wen sind neu:

- Automaten
- reguläre Ausdrücke
- formale Sprachen
(einschließlich reguläre Sprachen)

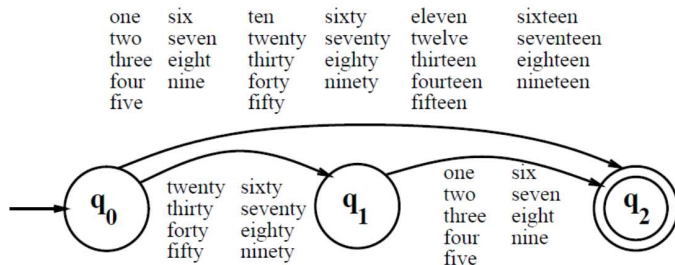
Übung

Denken Sie sich eine Endung des Adjektivs “klein” aus und finden Sie den Pfad, mit dem der Automat diese Endung erkennt/akzeptiert. Epsilon = “Nullübergang”. Endzustand: Doppelkreis.



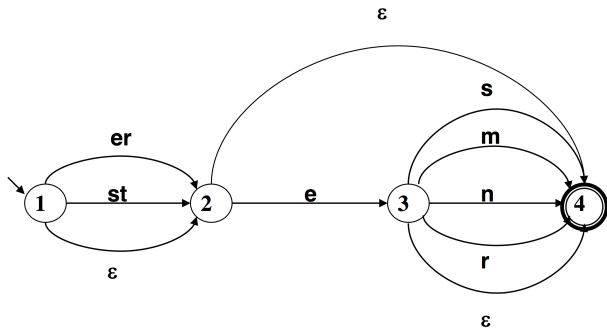
- Bogen, der mit einem Symbol annotiert ist:
Der Automat liest das Symbol und geht zum nächsten Zustand.
- Bogen, der mit einem Epsilon annotiert ist:
Der Automat liest nichts und geht zum nächsten Zustand.
- Doppelkreis: Endzustand.
Der Automat “akzeptiert” einen String, wenn er nach Lesen des Strings in einem Endzustand ist.
- Ein einfacher Pfeil (der nicht bei einem Zustand beginnt) zeigt auf den Startzustand.

Beispiel



Übung

Denken Sie sich eine Endung des Adjektivs “klein” aus und finden Sie den Pfad, mit dem der Automat diese Endung erkennt/akzeptiert. Epsilon = “Nullübergang”. Endzustand: Doppelkreis.



- Reguläre Sprachen
- Reguläre Ausdrücke
- deterministische Automaten
- nichtdeterministische Automaten
- Reguläre Sprachen
 - = Reguläre Ausdrücke
 - = deterministische Automaten
 - = nichtdeterministische Automaten

Caveat

In dieser Vorlesung gibt es einiges, das Sie als Hintergrund wissen sollten, das aber nicht abgefragt wird. Versuchen Sie einfach, möglichst viel zu verstehen, aber vollständiges Verständnis ist bei vielem nicht nötig.

- 1 Motivation
- 2 Reguläre Sprachen**
- 3 Reguläre Ausdrücke
- 4 Automaten
- 5 Breadth-First & Depth-First
- 6 Moodle

Das **Alphabet** Σ ist eine endliche Menge, deren Elemente **Symbole** oder **Zeichen** heißen. Wir verwenden für Zeichen üblicherweise die Buchstaben $a, b, c, \dots \in \Sigma$. Wir können die Zeichen aus einem gegebenen Alphabet Σ zu endlichen **Symbolfolgen** oder **Wörtern** zusammensetzen, wobei Wiederholungen erlaubt sind. Wir bezeichnen Wörter mit u, v, w, \dots oder griechischen Buchstaben. Die **Länge** $|w|$ eines Wortes w ist die Anzahl der Zeichen, aus denen es besteht. Eine besondere Stellung nimmt dabei das **leere Wort** ϵ ein, das als einziges die Länge $|\epsilon| = 0$ hat. Die **Menge aller Wörter** über einem Alphabet Σ (einschließlich des leeren Wortes) wird mit Σ^* bezeichnet.

Buchstabenebene

Zeichen sind Buchstaben.

Ein Wort der regulären Sprache ist ein **Wort** der natürlichen Sprache.

Lexemformebene

Zeichen sind Lexemformen (natürlichsprachliche Wörter).

Ein Wort der regulären Sprache ist eine **Phrase** oder ein **Satz** der natürlichen Sprache.

Formal läßt sich ein Wort $w \in \Sigma^*$ der Länge $|w| = n$ als Abbildung $w : \{1, \dots, n\} \rightarrow \Sigma$ auffassen.

Wir verwenden die Notation $w(i)$ für das i -te Zeichen eines Wortes: $w(1)$ bezeichnet also das erste, $w(n)$ das letzte Zeichen von w .

Alternativ: w_i, w_1, w_n

Verkettung = Konkatenation

Die **Verkettung** oder **Konkatenation** zweier Wörter $v, w \in \Sigma^*$ wird $v \circ w$ oder kurz vw geschrieben. $v \circ w \in \Sigma^*$ entsteht durch “Aneinanderhängen” der Wörter v und w . Wir können jedes Wort als Verkettung seiner Zeichen auffassen:
 $w = w(1) \circ w(2) \circ \dots \circ w(n)$ (mit $n = |w|$). Das leere Wort ϵ wird als **neutrales Element** (bezüglich der Konkatenation) bezeichnet, da $\epsilon \circ w = w \circ \epsilon = w$ für alle $w \in \Sigma^*$ gilt. Weiterhin ist stets $|v \circ w| = |v| + |w|$.

Die i -te **Potenz** w^i eines Wortes w bezeichnet dessen i -fache Wiederholung. Formal läßt sich dieser Sachverhalt in eine rekursive Definition fassen: $w^0 := \epsilon$ und $w^{i+1} := w^i \circ w$ für $i \in \mathbb{N}_0$. Die **Umkehrung** (oder **Spiegelung**) w^R von $w = w(1)w(2) \cdots w(n)$ ist definiert als $w^R := w(n)w(n-1) \cdots w(1)$ (z.B. ist $(\text{stefan})^R = \text{nafets}$). Ein **Palindrom** ist ein Wort w , für das $w^R = w$ gilt. Einfache Eigenschaften der Umkehrung sind $(v \circ w)^R = w^R \circ v^R$ und $|w^R| = |w|$.

Gibt es zu $v, w \in \Sigma^*$ ein Wort $u \in \Sigma^*$ mit $w = v \circ u$ (bzw. $w = u \circ v$), so heißt v **Präfix** (bzw. **Suffix**) von w . Ein Präfix v von w wird abkürzend mit $v \leq w$ gekennzeichnet.

Der **Abschluß** A^* (oder Abschluß unter Konkatenation) einer Wortmenge $A \subseteq \Sigma^*$ ist die kleinste Teilmenge von Σ^* , die folgende Bedingungen erfüllt:

- $A \subseteq A^*$
- $\epsilon \in A^*$
- $v, w \in A^* \Rightarrow v \circ w \in A^*$

Der Abschluß einer Wortmenge wird als **Kleenesche Hülle** bezeichnet, der Operator $*$ als **Kleene Star**. Anschaulich besteht A^* aus allen möglichen Verkettungen von beliebig vielen Wörtern aus A . Es gilt folglich

$$A^* = \{w_1 \circ w_2 \circ \dots \circ w_n \mid n \in \mathbb{N} \text{ und } w_1, \dots, w_n \in A\} \cup \{\epsilon\}$$

Definition: Formale Sprache

Eine **formale Sprache** über dem Alphabet Σ ist eine (beliebige) Teilmenge $A \subseteq \Sigma^*$. Wir unterscheiden zwischen **endlichen** und **unendlichen** Sprachen, je nachdem ob $|\Sigma|$ endlich oder unendlich ist. Da das Alphabet Σ endlich ist, ist Σ^* und somit jede unendliche formale Sprache abzählbar unendlich.

Auf formale Sprachen $A, B \subseteq \Sigma^*$ lassen sich die folgenden Mengenoperationen anwenden: **Vereinigung** $A \cup B$, **Durchschnitt** $A \cap B$, **Differenz** $A \setminus B$ sowie das **Komplement** $\mathcal{C}A = \Sigma^* \setminus A$ sind wieder formale Sprachen über Σ . Zusätzliche Operationen für formale Sprachen sind **Kleenesche Hülle** A^* und die **Verkettung** $A \circ B := \{v \circ w \mid v \in A \wedge w \in B\}$. Formal ist $A \circ B$ als Teilmenge von Σ^* definiert:

$$A \circ B := \{w \in \Sigma^* \mid \exists u \in A \exists v \in B : w = u \circ v\}.$$

Eine **Sprachklasse** \mathcal{L} über einem Alphabet Σ ist eine Menge von formalen Sprachen $A \subseteq \Sigma^*$. Da jede Sprache A selbst eine Teilmenge von Σ^* ist, ist \mathcal{L} also ein Menge von Teilmengen, d.h. eine Teilmenge der Potenzmenge von Σ^* : $\mathcal{L} \subseteq \mathcal{P}(\Sigma^*)$. Auch $\mathcal{P}(\Sigma^*)$ stellt somit eine Sprachklasse dar; sie umfaßt alle möglichen formalen Sprachen über Σ . Wir stellen an jede betrachtete Sprachklasse \mathcal{L} bestimmte Forderungen. Sie sind zugleich Minimalforderungen an die Mächtigkeit der entsprechenden Grammatikformalismen.

Hier: Die Sprachklasse der regulären Sprachen

Forderungen an eine Sprachklasse \mathcal{L}

- (i) Jede endliche Sprache gehört zu \mathcal{L} .
- Mit zwei Sprachen $A, B \in \mathcal{L}$ gehört auch deren (ii) Vereinigung $A \cup B$ zu \mathcal{L} und (iii) Verkettung $A \circ B$ zu \mathcal{L} .
- (iv) Mit jeder Sprache $A \in \mathcal{L}$ gehört auch deren Kleenesche Hülle A^* zu \mathcal{L} .

Forderung (i) stellt sicher, dass jeder Grammatikformalismus extensionale Definitionen ermöglicht. Forderungen (ii) und (iii) leiten sich aus der Kombination der Regeln zweier Grammatiken her, wobei die Regelsätze entweder Alternativen sind (ii) oder durch Verkettung kombiniert werden (iii). Forderung (iv) schließlich bildet den Schlüssel zu unendlichen Sprachen. Während (iii) bereits sicherstellt, dass Verkettungen endlicher Länge $\in \mathcal{L}$ sind, wird diese Eigenschaft durch (iv) auf Verkettungen beliebiger Länge ausgedehnt.

Die einfachste Sprachklasse, die obigen Forderungen genügt, bilden die sogenannten **regulären Sprachen**. Wir bezeichnen die Klasse aller regulären Sprachen über einem gegebenen Alphabet Σ mit $\text{Reg}(\Sigma) \subseteq \mathcal{P}(\Sigma^*)$.

Wir können $\text{Reg}(\Sigma)$ analog zu den Forderungen (i) – (v) rekursiv definieren:

- (i) $\emptyset, \{\epsilon\} \in \text{Reg}(\Sigma)$
- (ii) $\forall a \in \Sigma : \{a\} \in \text{Reg}(\Sigma)$
- (iii) $A, B \in \text{Reg}(\Sigma) \Rightarrow A \cup B \in \text{Reg}(\Sigma)$
- (iv) $A, B \in \text{Reg}(\Sigma) \Rightarrow A \circ B \in \text{Reg}(\Sigma)$
- (v) $A \in \text{Reg}(\Sigma) \Rightarrow A^* \in \text{Reg}(\Sigma)$

Forderung (i) wird durch die Regeln (i), (ii) und (iii) erfüllt, da mit ihnen jede endliche Sprache erzeugt werden kann. Jede reguläre Sprache $A \in \text{Reg}(\Sigma)$ läßt sich durch wiederholte Anwendung der obigen Regeln darstellen.

Die Sprache $A := \{ab^k c^l \mid k \in \mathbb{N}_0, l \in \{0, 1\}\}$ ist regulär, denn

$$A = \{a\} \circ (\{b\})^* \circ (\{c\} \cup \{\epsilon\}).$$

Das Hauptinteresse der **Theorie formaler Sprachen** besteht darin, eine **unendliche** Sprache durch eine **endliche Grammatik** zu beschreiben. Hierzu werden geeignete **Beschreibungsformalismen** benötigt. Die Gesamtheit aller Sprachen A , die sich in einem bestimmten Formalismus beschreiben lassen, wird als **Sprachklasse** \mathcal{A} bezeichnet. Ein wesentlicher Aspekt der Forschungsarbeit ist somit die Untersuchung der Eigenschaften verschiedener Sprachklassen.

- Reguläre Sprachen
- Reguläre Ausdrücke
- deterministische Automaten
- nichtdeterministische Automaten
- Reguläre Sprachen
 - = Reguläre Ausdrücke
 - = deterministische Automaten
 - = nichtdeterministische Automaten

- 1 Motivation
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke**
- 4 Automaten
- 5 Breadth-First & Depth-First
- 6 Moodle

Reguläre Ausdrücke (1)

Eine etwas kompaktere Darstellung von regulären Sprachen als die gerade eingeführte Mengenschreibweise erlauben **reguläre Ausdrücke**. Ein regulärer Ausdruck über einem Alphabet Σ besteht aus den Zeichen von Σ sowie den sogenannten **Metazeichen** $(,), |$ und $*$. Um Verwechslungen von regulären Ausdrücken mit Wörtern über Σ zu vermeiden, werden reguläre Ausdrücke auf diesen Folien durch **Unterstreichung** gekennzeichnet (z.B. ist aab ein Wort, aab ein regulärer Ausdruck). Wir verwenden r , s , t , ... als Variablen für reguläre Ausdrücke.

Reguläre Ausdrücke (2)

Wir können reguläre Ausdrücke über Σ als Wörter über dem Alphabet $\Sigma \cup M$ auffassen, wobei $M := \{ (,), |, * \}$ die Menge der Metazeichen ist und $\Sigma \cap M = \emptyset$ angenommen wird. Die Menge aller zulässigen regulären Ausdrücke über Σ bezeichnen wir mit $R(\Sigma)$. Jeder Ausdruck $\underline{r} \in R(\Sigma)$ beschreibt eine formale Sprache über Σ , die wir mit $\mathcal{L}[\underline{r}]$ bezeichnen.

Abbildung: Regulärer Ausdruck \mapsto Reguläre Sprache

Die Menge aller regulären Ausdrücke $R(\Sigma)$ und die **Abbildung** $\mathcal{L} : R(\Sigma) \rightarrow \mathcal{P}(\Sigma^*)$, die jedem regulären Ausdruck die durch ihn beschriebene Sprache zuordnet, werden wiederum rekursiv definiert:

$\epsilon \in R(\Sigma)$	$\mathcal{L}[\underline{\epsilon}] := \{\epsilon\}$	
$\forall a \in \Sigma : \underline{a} \in R(\Sigma)$	$\mathcal{L}[\underline{a}] := \{a\}$	
$\underline{r}, \underline{s} \in R(\Sigma) \Rightarrow \underline{(r)(s)} \in R(\Sigma)$	$\mathcal{L}[\underline{(r)(s)}] := \mathcal{L}[\underline{r}] \circ \mathcal{L}[\underline{s}]$	Konkatenation
$\underline{r}, \underline{s} \in R(\Sigma) \Rightarrow \underline{r s} \in R(\Sigma)$	$\mathcal{L}[\underline{r s}] := \mathcal{L}[\underline{r}] \cup \mathcal{L}[\underline{s}]$	Disjunktion
$\underline{r} \in R(\Sigma) \Rightarrow \underline{(r)*} \in R(\Sigma)$	$\mathcal{L}[\underline{(r)*}] := (\mathcal{L}[\underline{r}])^*$	Kleene-Hülle

Aus dieser Definition ist unmittelbar ersichtlich, dass reguläre Ausdrücke gerade die regulären Sprachen beschreiben, bzw. genauer:

Satz

$$\text{Reg}(\Sigma) = \mathcal{L}[R(\Sigma)] \cup \{\emptyset\}.$$

Zur Vereinfachung der Notation vereinbaren wir, dass Klammern weggelassen werden dürfen, wo sie aufgrund der Assoziativität von Disjunktion und Konkatenation bzw. aufgrund der Präzedenzregeln (Kleenesche Hülle vor Konkatenation vor Disjunktion) nicht benötigt werden (ebenso dürfen zusätzliche Klammern eingefügt werden, um die Lesbarkeit zu verbessern). Beispielsweise ist wegen der Assoziativität der Konkatenation

$$\underline{((a)(b))(c)} \equiv \underline{(a)((b)(c))} \equiv: \underline{abc}$$

Zwei reguläre Ausdrücke $\underline{r}, \underline{s} \in R(\Sigma)$ sind **äquivalent** ($\underline{r} \equiv \underline{s}$), wenn $\mathcal{L}[\underline{r}] = \mathcal{L}[\underline{s}]$.

Es ist nicht immer offensichtlich, ob zwei reguläre Ausdrücke äquivalent sind, z.B.

$$\underline{(ab^*)^*} \equiv \underline{(a(a^*b^*)^*)} | \epsilon \equiv \underline{(a(a|b)^*)} | \epsilon$$

Daher ist es wünschenswert, reguläre Ausdrücke auf eine sogenannte “Normalform” zurückzuführen. Jeder zu einem regulären Ausdruck $\underline{r} \in R(\Sigma)$ äquivalente Ausdruck der Form $\underline{s_1 | s_2 | \dots | s_n}$, wobei keiner der Ausdrücke $\underline{s_i}$ das Metazeichen $|$ enthält, heißt **disjunktive Normalform** von \underline{r} . Die disjunktive Normalform ist im allgemeinen **nicht eindeutig** bestimmt.

Für Software-Anwendungen erweist sich die bisher verwendete **Standardnotation** regulärer Ausdrücke oft als unhandlich. Sie wurde daher nach und nach um weitere Metazeichen ergänzt. Inzwischen hat sich als Standard die sogenannte **POSIX-Notation** mit den Metazeichen

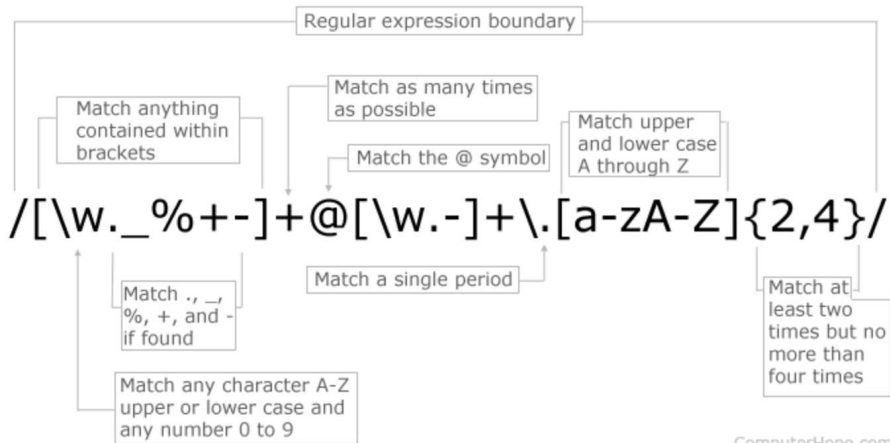
$$M := \{ (,), |, *, ., ?, +, \{, \}, [, -,], ^, \$ \dots \}$$

etabliert.

Das Alphabet Σ besteht stets aus dem kompletten Zeichenvorrat der Anwendung (ASCII, Unicode, ...), daher ist zwangsläufig $M \cap \Sigma \neq \emptyset$ (und sogar $M \subseteq \Sigma$). Soll ein Zeichen aus M "wörtlich", d.h. als Zeichen von Σ interpretiert werden, so wird ein **Backslash** \backslash vorangestellt ($\$, \wedge, ,$ und $-$ werden nur in bestimmten Kontexten als Metazeichen interpretiert).

Man beachte, dass der leere reguläre Ausdruck ϵ nicht durch ein spezielles Zeichen sondern durch einen leeren String repräsentiert wird. So würde beispielsweise der POSIX-Ausdruck $[ac] + (baa|)$ in der Standardnotation $(a|c)(a|c)^*(baa|\epsilon)$ lauten.

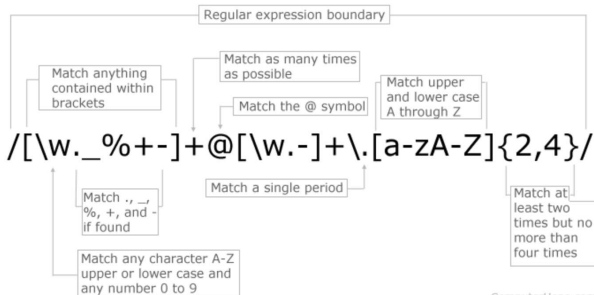
Regulärer Ausdruck für Email-Adressen



ComputerHope.com

Übung

- (i) Modifiziere den Ausdruck, so dass die neuen top level domains richtig erkannt werden. (ii) Denken Sie sich eine valide Email-Adresse aus, die vom neuen Ausdruck **nicht** erkannt wird. (iii) Gibt es eine nicht valide Email-Adresse, die vom neuen Ausdruck trotzdem als richtig erkannt wird?



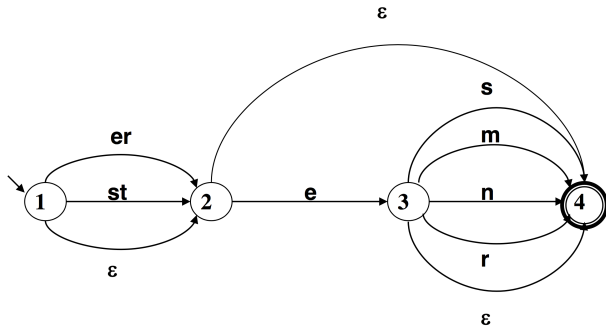
ComputerHope.com

- Reguläre Sprachen
- Reguläre Ausdrücke
- deterministische Automaten
- nichtdeterministische Automaten
- Reguläre Sprachen
 - = Reguläre Ausdrücke
 - = deterministische Automaten
 - = nichtdeterministische Automaten

- 1 Motivation
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Automaten**
- 5 Breadth-First & Depth-First
- 6 Moodle

Ein **endlicher Automat (EA)** M über dem Alphabet Σ ist formal ein Quintupel $M = (Q, \Sigma, s, F, \delta)$. Dabei bezeichnet Q eine endliche Menge von **Zuständen**, $s \in Q$ den **Startzustand** des Automaten und $F \subseteq Q$ die Menge der **Endzustände**. Die **Übergangsrelation** $\delta \subseteq Q \times \Sigma \times Q$ beschreibt, welche Zustandsübergänge der EA ausführen kann: ist $(q, a, q') \in \delta$, so kann M durch Einlesen des Zeichens a (**Label**) vom Zustand q (**Ausgangszustand**) in den Zustand q' (**Zielzustand**) wechseln. Wir schreiben Übergangsregeln in der Form $q \xrightarrow{a} q'$ und bezeichnen sie kurz als **Bewegungen**.

Automat für Adjektive (Buchstabenebene)



Gibt es für jeden Ausgangszustand q und jedes Zeichen a höchstens einen Zielzustand (d.h. höchstens eine Übergangsregel $q \xrightarrow{a} q'$), dann bezeichnen wir M als **deterministisch (DEA)**, andernfalls als **nichtdeterministisch (NEA)**. Für einen deterministischen EA können wir die Übergangsrelation als (partielle) Abbildung $\delta : Q \times \Sigma \rightarrow Q$ auffassen.

Bei einem NEA erlauben wir Epsilon-Bewegungen, also Zustandsübergänge ohne Einlesen eines Zeichens. Diese werden als Bewegungen mit dem Label ϵ dargestellt ($q \xrightarrow{\epsilon} q'$), also ist hier $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$.

Beispiel

Deterministischer EA $M_1 := (Q_1, \Sigma, s_1, F_1, \delta_1)$ über $\Sigma = \{a, b\}$ mit

$$Q_1 := \{q_0, q_1, q_2\}$$

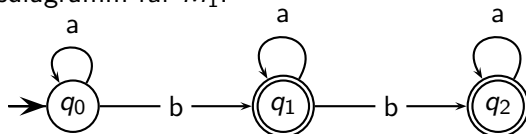
$$s_1 := q_0$$

$$F_1 := \{q_1, q_2\}$$

$$\delta_1 := \{q_0 \xrightarrow{a} q_0, q_0 \xrightarrow{b} q_1, q_1 \xrightarrow{a} q_1, q_1 \xrightarrow{b} q_2, q_2 \xrightarrow{a} q_2\}.$$

Endliche Automaten lassen sich anschaulich durch

Übergangsdiagramme darstellen. Die Abbildung zeigt das Übergangsdiagramm für M_1 .



Ein Automat akzeptiert ein Eingabewort, wenn . . .

Ein Wort $w \in \Sigma^*$ wird von dem EA M **akzeptiert**, wenn dieser sich nach Einlesen aller Zeichen von w in einem Endzustand $q \in F$ befindet. Umgekehrt wird w **nicht akzeptiert**, wenn M sich zum Schluß nicht in einem Endzustand befindet oder schon während der Verarbeitung von w mit einem Fehler abbricht, weil keine zu dem nächsten einzulesenden Zeichen passende Bewegung definiert ist. Gewöhnlich wird kein Unterschied zwischen den beiden Arten des Nichtakzeptierens gemacht.

Die schrittweise Abarbeitung eines Eingabewortes w auf dem Rechner nennt man **Simulation** des Automaten M . Der Zustand, in dem sich M jeweils befindet, wird als **aktiver Zustand** bezeichnet. Bei einem nichtdeterministischen EA können mehrere Zustände gleichzeitig aktiv sein. Die Menge aller akzeptierten Wörter $w \in \Sigma^*$ bezeichnet man als die **von M akzeptierte Sprache** $\mathcal{L}[M]$.

Zur formalen Definition von $\mathcal{L}[M]$ führen wir den Begriff der **Konfiguration** $(q, w) \in Q \times \Sigma^*$ ein, wobei q den gerade aktiven Zustand und w die noch zu verarbeitende Eingabe bezeichnet. Der **Übergang** $(q, a \circ w) \vdash (q', w)$ von einer Konfiguration $(q, a \circ w)$ in die Konfiguration (q', w) (d.h. der Wechsel vom Zustand q in den Zustand q' unter Einlesen des Zeichens a) ist **erlaubt**, wenn eine entsprechende Bewegung $q \xrightarrow{a} q' \in \delta$ definiert ist. Im Falle eines ϵ -EA erlaubt jede ϵ -Bewegung $q \xrightarrow{\epsilon} q'$ einen Übergang $(q, w) \vdash (q', w)$ ohne Einlesen eines Zeichens.

Ein **Pfad** von (q, w) nach (q', w') ist eine Kette von Übergängen

$$(q, w) = (q_0, w_0) \vdash (q_1, w_1) \vdash \cdots \vdash (q_n, w_n) = (q', w').$$

Wir schreiben kurz $(q, w) \vdash^* (q', w')$, wenn (mindestens) ein solcher Pfad existiert.

Wenden wir M auf das Eingabewort w an, so startet der Automat in der Konfiguration (s, w) und w wird genau dann akzeptiert, wenn es einen Pfad zu einem geeigneten Endzustand $q \in F$ gibt, der w vollständig abarbeitet: $(s, w) \vdash^* (q, \epsilon)$. Also können wir definieren:

$$\mathcal{L}[M] := \{w \in \Sigma^* \mid \exists q \in F : (s, w) \vdash^* (q, \epsilon)\}$$

Zwei EA M_1 und M_2 über demselben Alphabet Σ heißen **äquivalent** ($M_1 \equiv M_2$), wenn $\mathcal{L}[M_1] = \mathcal{L}[M_2]$.

Es ist unmittelbar klar, dass die DEA Spezialfälle der NEA sind. Somit gibt es zu jedem DEA M einen äquivalenten NEA. Um die Äquivalenz von NEA und DEA nachzuweisen, müssen wir also noch zeigen, dass es zu jedem NEA einen äquivalenten DEA.

Lemma

Sei $M = (Q, \Sigma, s, F, \delta)$ ein NEA. Dann gibt es einen zu M äquivalenten DEA $M' = (Q', \Sigma, s', F', \delta')$, d.h. es gilt $\mathcal{L}[M'] = \mathcal{L}[M]$.

Equivalence of NFAs and DFAs

- The ϵ -closure of a state q , denoted $E(q)$, is the set of all states, including q itself, that can be reached using only ϵ -moves.
- Let $\{Q, \Sigma, \delta, q_0, F\}$ be an NFA.
- Define a DFA $\{Q', \Sigma', \delta', q'_0, F'\}$ as follows:
 - $Q' = \mathcal{P}(Q)$.
 - $\Sigma' = \Sigma$.
 - $q'_0 = E(q_0)$.
 - For any state $S \in \mathcal{P}(Q)$ and any $a \in \Sigma$, define $\delta'(S, a)$ to be

$$\delta'(S, a) = \bigcup_{q \in S} E(\delta(q, a)).$$

- The final states in F' are those states that contain final states of the NFA.

“Automatensprachen” = reguläre Sprachen

Endliche Automaten akzeptieren genau diejenigen Sprachen, die durch reguläre Ausdrücke beschrieben werden können, also die Klasse der **regulären Sprachen**. Zum Nachweis dieser Behauptung müssen wir die **Äquivalenz** von EA und regulären Ausdrücken zeigen: zu jedem regulären Ausdruck \underline{r} gibt es (mindestens) einen endlichen Automaten M mit $\mathcal{L}[M] = \mathcal{L}[\underline{r}]$; umgekehrt gibt es zu jedem endlichen Automaten M einen regulären Ausdruck \underline{r} mit $\mathcal{L}[\underline{r}] = \mathcal{L}[M]$.

Zu jedem regulären Ausdruck \underline{r} gibt es (mindestens) einen endlichen Automaten M mit $\mathcal{L}[M] = \mathcal{L}[\underline{r}]$:

Konstruktionsverfahren.

Zu jedem endlichen Automaten M gibt es einen regulären Ausdruck \underline{r} mit $\mathcal{L}[\underline{r}] = \mathcal{L}[M]$. Diese Richtung ist komplexer und sehr “technisch”, d.h. intuitiv schwer nachvollziehbar.

- Reguläre Sprachen
- Reguläre Ausdrücke
- deterministische Automaten
- nichtdeterministische Automaten
- Reguläre Sprachen
 - = Reguläre Ausdrücke
 - = deterministische Automaten
 - = nichtdeterministische Automaten

- 1 Motivation
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Automaten
- 5 Breadth-First & Depth-First**
- 6 Moodle

Breadth first search

- Breadth first traversal is accomplished by enqueueing each level of a tree sequentially as the root of any subtree is encountered. There are 2 cases in the iterative algorithm.
- Root case: The traversal queue is initially empty so the root node must be added before the general case.
- General case: Process any items in the queue, while also expanding their children. Stop if the queue is empty. The general case will halt after processing the bottom level as leaf nodes have no children.
- Input: A search problem. A search-problem abstracts out the problem specific requirements from the actual search algorithm.
- Output: An ordered list of actions to be followed to reach from start state to the goal state.

Wikipedia

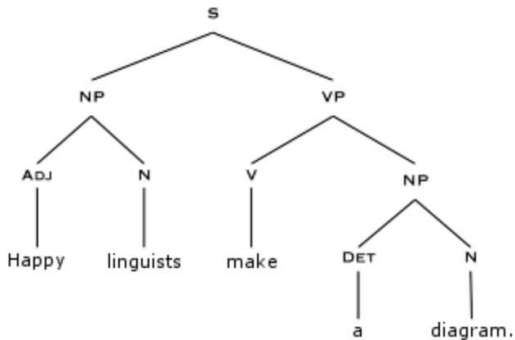
Depth first search (DFS)

```
procedure DFS(G,v):  
  process v  
  label v as discovered  
  for all edges from v to w in G.adjacentEdges(v) do  
    if vertex w is not labeled as discovered then  
      recursively call DFS(G,w)
```

Wikipedia

Übung

Give order of nodes for BFS and DFS



- Reguläre Ausdrücke (POSIX)
- Formale Definition des DEA
- Formale Definition des NEA
- Unterschiede DEA vs NEA
- Automaten zeichnen
- Depth first search (DFS)
- Breadth first search (BFS)
- Nicht relevant: Determinisierung von NEAs

- 1 Motivation
- 2 Reguläre Sprachen
- 3 Reguläre Ausdrücke
- 4 Automaten
- 5 Breadth-First & Depth-First
- 6 Moodle**

In der folgenden Aufgabe werden Sie mit dem Tiger-Corpus arbeiten. Das Tiger-Corpus enthält ca. 50000 deutsche Sätze, die mit morphologischen und syntaktischen Informationen wie z.B. Wortarten annotiert sind. Auf Moodle ist ein Ausschnitt aus dem Tiger-Corpus verlinkt ([tiger.txt](#)), in dem in jeder Zeile ein Wort steht.

.	irgendein beliebiges Zeichen
.*	kein oder beliebig viele Zeichen
[a-e]	a, b, c, d oder e
[^a-e]	alle Zeichen außer a, b, c, d und e
(maus hund)	Zeichenfolge maus oder hund
(ab)*	kein oder beliebig viele ab
(ab)+	mindestens ein oder oder mehr ab
(ab)?	kein oder ein ab
^	Anfang einer Zeile
\$	Ende einer Zeile

Unter Linux können Sie mit dem Kommando `grep` mit regulären Ausdrücken in Dateien suchen (unter Windows geht das mit Notepad++ in der erweiterten Suchoption). Das Kommando `grep -P "abc" datei.txt` gibt z.B. jede Zeile von `datei.txt` aus, in der `abc` vorkommt.

Versuchen Sie, einen regulären Ausdruck zu finden, mit dem Sie alle Wortformen (Präsens, Präteritum, mit Partizipialendungen) des Wortes reden (aber nicht z.B. von mitreden oder Gerede) im Tiger-Corpus finden, und geben Sie ihn an. Unter Linux können Sie außerdem wc verwenden, um herauszufinden, wie viele Vorkommen es sind.

Einige Sonderzeichen (z.B. wenn nach "+" gesucht werden soll) müssen mit einem Backslash (nach links geneigter Schrägstrich: "\") escaped werden; dazu und generell zu regulären Ausdrücken gibt es viele vollständige und ausführliche Hilfeseiten online.

Spielen Sie ein wenig mit `grep` und anderen Unix-Befehlen! der `man`-Befehl (z.B. `man grep`) liefert zu fast allen Unix-Befehlen Erklärungen und eine Auflistung möglicher Optionen. Übung macht den Meister, und mit ein wenig Tüftelei kann man in einer Zeile mächtige Befehle zusammenstellen.

Mit dem Operator `|` können mehrere Unix-Befehle hintereinandergeschaltet werden (man spricht auch von "pipen" oder einer "Pipeline"), sodass das Ergebnis des einen Befehls vom nächsten Befehl aufgenommen und weiterverarbeitet wird.

Mit den Operatoren `>` und `>>` kann man das Ergebnis eines Befehls in eine Datei schreiben bzw. an eine Datei anhängen, um nicht alles auf der Konsole ausgegeben zu bekommen.

In regulären Ausdrücken werden (runde Klammern) für zusammengehörige Zeichenfolgen (= "Gruppen") verwendet und [eckige Klammern], um Mengen/Klassen von Zeichen zu deklarieren. Gruppen sind innerhalb von Klassen nicht aussagekräftig. Dies kann zu Verwirrung führen. Um mehrere Gruppen zu quantifizieren (z.B. mit +), kann man diese zu einer "veroderten Gruppe" machen: $((\text{gruppe1})|(\text{gruppe2})|(\text{gruppe3}))+$ z.B. matcht alles, was mindestens eine der drei Gruppen enthält, egal in welcher Reihenfolge oder Zusammensetzung.

Um sicherzustellen, dass alle Wortformen von reden gematcht werden, sollten Sie auch wirklich alle Wortformen testen und nicht darauf vertrauen, dass wirklich alle Wortformen von reden im Tiger-Corpus vorkommen.

```
grep 'red' tiger.txt | uniq | sort | uniq -c | sort -gr > output.txt
```