

# 有限状态机与最小编辑距离

## 01 有限状态机

### 有限状态机

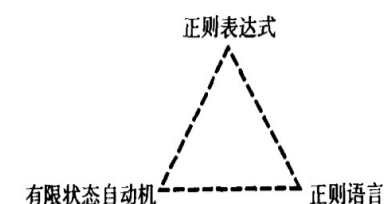
1. 概述

2. 应用

stefanie

### 有限状态机的概述

- 任何正则表达式都可以用有限状态自动机来实现(除了使用存储特性的那些正则表达式之外)
- 任何有限状态自动机都可以用正则表达式来描述。有限状态自动机和正则表达式彼此对称。
- 正则表达式是用来刻画正则语言 (regular language) 的一种方法
- 正则语言是一种特别的形式语言。



stefanie

## 有限状态机的概述

### 如何识别羊的叫声

- baa!
- baaa!
- baaaa!
- baaaaa!
- baaaaaa!..

问题：正则表达式？

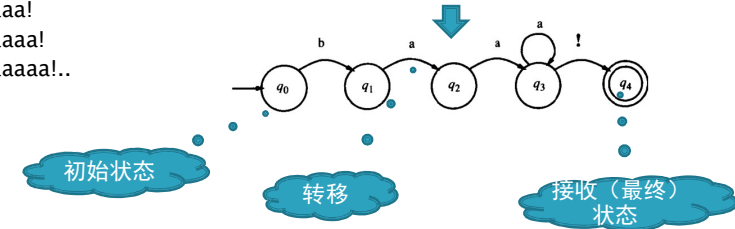
stefanie

## 有限状态机的概述

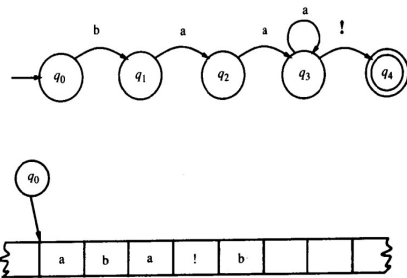
### 如何识别羊的叫声

- baa!
- baaa!
- baaaa!
- baaaaa!
- baaaaaa!..

**有限状态机** (finite-state machine, 缩写: **FSM**) 又称有限状态自动机, 简称状态机, 是表示有限个状态以及在这些状态之间的转移和动作等行为的数学模型。

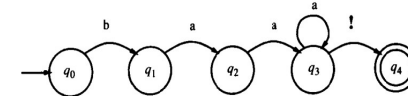


## 有限状态机的概述



- 从形式上说, 一个有限自动机可以用下面5个参数来定义:
- **Q**: N中状态  $q_0, \dots, q_w$  的有限集合
  - **$\Sigma$** : 有限的输入符号字母表
  - **$q_0$** : 初始状态
  - **F**: 终极状态的集合,  $F \subseteq Q$
  - **$\delta(q, i)$** : 状态之间的转移函数或转移矩阵。给定一个状态  $q \in Q$  和一个输入符号  $i \in \Sigma$ ,  $\delta(q, i)$  返回一个新的状态  $q'$ , 因此,  $\delta(q, i)$  是从  $Q \times \Sigma$  到  $Q$  的一个关系。

## 有限状态机的概述



- 上例中,
- $Q = ?$
- $\Sigma = ?$ ,
- $F = ?$

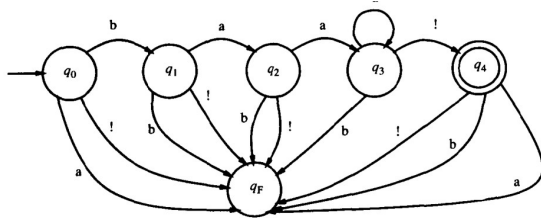
**$\delta(q, i)$  用转移表确定:**

状态	输出		
	b	a	!
0	1	0	0
1	0	2	0
2	0	3	0
3	0	3	4
4:	0	0	0

stefanie

## 有限状态机的概述

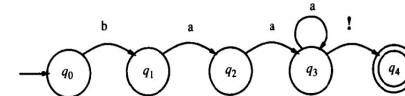
- 增加一个失败状态：



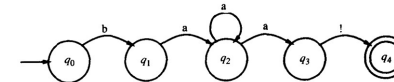
stefanie

## 有限状态机的概述

- 确定性的自动机 (Deterministic FSA或DFAS),



- 非确定的有限自动机 (non-deterministic FSA或NFA)



stefanie

## 有限状态机的应用

- 地址识别

广东省深圳市腾讯大厦

广东省 518057 深圳市南山区科技园腾讯大厦

深圳市 518057 科技园腾讯大厦

深圳市南山区科技园腾讯公司

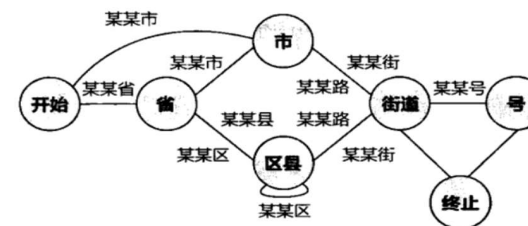
深圳市南山区科技园腾讯总部 518000 (估计不知道准确的邮编)

广东省深圳市科技园中一路腾讯公司

.....

[https://blog.csdn.net/qq\\_16234613](https://blog.csdn.net/qq_16234613)

## 有限状态机的应用

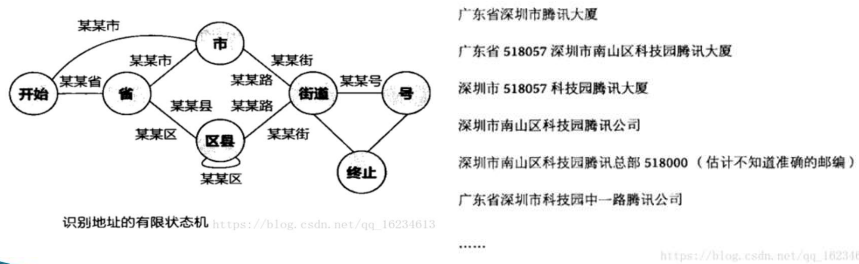


识别地址的有限状态机

如果一个地址能从状态机的起始状态，经过若干中间状态，走到最终状态，那么这条地址就是有效的

stefanie

- 思考：根据有限转态机右下图的哪些地址可以识别？



stefanie

## 有限状态机的应用

- 有限自动机用于英语单词形态分析[Allen, 1995]

英语单词形态变化非常普遍，例如：

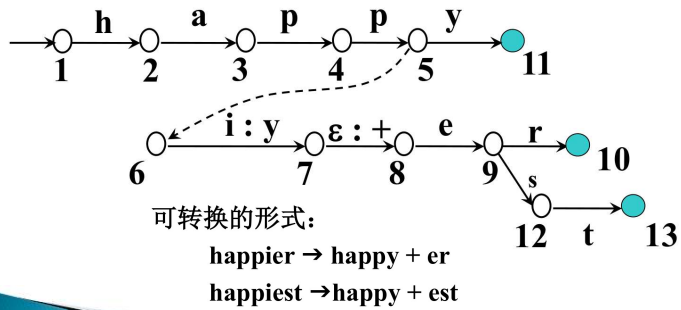
**eat:** eats, eating, ate, eaten

**happy:** happier, happiest

stefanie

## 有限状态机的应用

- 任务：识别单词:happy, happier, happiest



stefanie

02

最小编辑距离

## 最小编辑距离

1. 引言
2. 最小编辑距离
3. 最小编辑距离-动态规划
4. 其他应用

stefanie

## 引言

### ▶ 拼写纠错

涛宝 → 淘宝

happy → hoppy

在英文查询串中，有Spelling Error的查询串占到了总查询串的10%~15% (Cucerzan & Brill 2004)，而中文有输入错误的查询串保守估计也要占到所有查询串的5%~10%左右



纠错算法主要是列出语料库词典中与原查询词具有最小编辑距离的词条，将其作为候选词

## 最小编辑距离

### ▶ 最小编辑距离

- 将一个错误拼写的单词被纠正成正确单词的最小编辑次数，每一次编辑只能改变一个字母。
- 俄罗斯科学家 Vladimir Levenshtein 在1965年提出来的，所以编辑距离又称为 Levenshtein距离

acress → actress (插入操作) : 1  
acress → access (替换操作) : 2  
acress → cress (删除操作) : 1

stefanie

## 最小编辑距离

- ▶ 原始串 s o t
- ▶ 目标串 s t o p

s o t  
→ s t o t  
→ s t o p

编辑距离: 3

s o t  
→ s t t  
→ s t o  
→ s t o p

编辑距离: 5

stefanie

## 最小编辑距离的动态规划

$$\text{substituteCost} \begin{cases} = 0 & \text{if } \text{target}[i] = \text{source}[j] \\ = 2 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \text{insertCost} &= 1 \\ \text{deleteCost} &= 1 \end{aligned}$$

$$\begin{aligned} D(0, 0) &= 0 \\ D(i, 0) &= \text{insertCost} * i \\ D(0, j) &= \text{deleteCost} * j \end{aligned}$$

i, 目标字符串位置  
j, 原始字符串位置

$$D(i, j) = \min \begin{cases} D(i-1, j) + \text{insertCost}(\text{target}_i) \\ D(i-1, j-1) + \text{substituteCost}(\text{source}_j, \text{target}_i) \\ D(i, j-1) + \text{deleteCost}(\text{source}_j) \end{cases}$$

Termination:  $D(N, M)$  is distance

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

$n = \text{length}(\text{target})$   
 $m = \text{length}(\text{source})$   
 Create matrix  $d[n, m]$ ;

$i=0 \quad j=0$

$d[0,0] = 0;$   
 $d[0,1] = 1; \dots; d[0,m] = m;$   
 $d[1,0] = 1; \dots; d[n,0] = n;$

j ↑	3	t				
	2	o				
	1	s				
	0	#	s	t	o	p
	#	0	1	2	3	4
						i →

## 最小编辑距离的动态规划

- 原始串 s: s o t
- 目标串 t: s t o p

$i=1 \quad j=1$

$$d[1,1] = \min \begin{cases} d[0,1] + \text{insert}(t[1]) = 2 \\ d[0,0] + \text{substitute}(s[1], t[1]) = 0 \\ d[1,0] + \text{delete}(s[1]) = 2 \end{cases} = 0$$

j ↑	3	t				
	2	o				
	1	s	0			
	0	#	s	t	o	p
	#	0	1	2	3	4
						i →

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

$i=1 \quad j=2$

$$d[1,2] = \min \begin{cases} d[0,2] + \text{insert}(t[1]) = 3 \\ d[0,1] + \text{substitute}(s[2], t[1]) = 3 \\ d[1,1] + \text{delete}(s[2]) = 1 \end{cases} = 1$$

j ↑	3	t				
	2	o	1			
	1	s	0			
	0	#	s	t	o	p
	#	0	1	2	3	4
						i →

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

3	t	2				
2	o	1				
1	s	0				
0	#	s	t	o	p	
#	0	1	2	3	4	

i=1 j=3

$$d[1,3] = \min \begin{cases} d[0,3] + \text{insert}(t[1]) = 4 \\ d[0,2] + \text{substitute}(s[3], t[1]) = 4 \\ d[1,2] + \text{delete}(s[3]) = 2 \end{cases} = 2$$

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

3	t	2				
2	o	1				
1	s	0	1			
0	#	s	t	o	p	
#	0	1	2	3	4	

i=2 j=1

$$d[2,1] = \min \begin{cases} d[1,1] + \text{insert}(t[2]) = 1 \\ d[1,0] + \text{substitute}(s[1], t[2]) = 3 \\ d[2,0] + \text{delete}(s[1]) = 3 \end{cases} = 1$$

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

3	t	2				
2	o	1	2			
1	s	0	1			
0	#	s	t	o	p	
#	0	1	2	3	4	

i=2 j=2

$$d[2,2] = \min \begin{cases} d[1,2] + \text{insert}(t[2]) = 2 \\ d[1,1] + \text{substitute}(s[2], t[2]) = 2 \\ d[2,1] + \text{delete}(s[2]) = 2 \end{cases} = 2$$

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

3	t	2	1			
2	o	1	2			
1	s	0	1			
0	#	s	t	o	p	
#	0	1	2	3	4	

i=2 j=3

$$d[2,3] = \min \begin{cases} d[1,3] + \text{insert}(t[2]) = 3 \\ d[1,2] + \text{substitute}(s[3], t[2]) = 1 \\ d[2,2] + \text{delete}(s[3]) = 3 \end{cases} = 1$$

stefanie



## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

j	3	t	2	1		
2	o	1	2			
1	s	0	1	2		
0	#	s	t	o	p	
#	0	1	2	3	4	

i=3 j=1

$$d[3,1] = \min \begin{cases} d[2,1] + \text{insert}(t[3]) = 2 \\ d[2,0] + \text{substitute}(s[1], t[3]) = 4 \\ d[3,0] + \text{delete}(s[1]) = 4 \end{cases} = 2$$

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

j	3	t	2	1		
2	o	1	2	1		
1	s	0	1	2		
0	#	s	t	o	p	
#	0	1	2	3	4	

i=3 j=2

$$d[3,2] = \min \begin{cases} d[2,2] + \text{insert}(t[3]) = 3 \\ d[2,1] + \text{substitute}(s[2], t[3]) = 1 \\ d[3,1] + \text{delete}(s[2]) = 3 \end{cases} = 1$$

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

j	3	t	2	1	2	
2	o	1	2	1		
1	s	0	1	2		
0	#	s	t	o	p	
#	0	1	2	3	4	

i=3 j=3

$$d[3,3] = \min \begin{cases} d[2,3] + \text{insert}(t[3]) = 2 \\ d[2,2] + \text{substitute}(s[3], t[3]) = 4 \\ d[3,2] + \text{delete}(s[3]) = 2 \end{cases} = 2$$

stefanie

## 最小编辑距离的动态规划

- 原始串 s o t
- 目标串 s t o p

j	3	t	2	1	2	
2	o	1	2	1		
1	s	0	1	2	3	
0	#	s	t	o	p	
#	0	1	2	3	4	

i=4 j=1

$$d[4,1] = \min \begin{cases} d[3,1] + \text{insert}(t[4]) = 3 \\ d[3,0] + \text{substitute}(s[1], t[4]) = 5 \\ d[4,0] + \text{delete}(s[1]) = 5 \end{cases} = 3$$

stefanie



## 最小编辑距离的动态规划

- 原始串 `s o t`
- 目标串 `s t o p`

$$j \downarrow \begin{array}{|c|c|c|c|c|c|} \hline 3 & t & 2 & 1 & 2 & \\ \hline 2 & o & 1 & 2 & 1 & 2 \\ \hline 1 & s & 0 & 1 & 2 & 3 \\ \hline 0 & \# & s & t & o & p \\ \hline \# & 0 & 1 & 2 & 3 & 4 \\ \hline \end{array} i \rightarrow$$

$i=4 \quad j=2$

$$d[4,2] = \min \left\{ \begin{array}{l} d[3,2] + \text{insert}(t[4]) = 2 \\ d[3,1] + \text{substitute}(s[2], t[4]) = 4 \\ d[4,1] + \text{delete}(s[2]) = 4 \end{array} \right\} = 2$$

stefanie

## 最小编辑距离的动态规划

- 原始串 `s o t`
- 目标串 `s t o p`

$$j \downarrow \begin{array}{|c|c|c|c|c|c|} \hline 3 & t & 2 & 1 & 2 & 3 \\ \hline 2 & o & 1 & 2 & 1 & 2 \\ \hline 1 & s & 0 & 1 & 2 & 3 \\ \hline 0 & \# & s & t & o & p \\ \hline \# & 0 & 1 & 2 & 3 & 4 \\ \hline \end{array} i \rightarrow$$

$i=4 \quad j=3$

$$d[4,3] = \min \left\{ \begin{array}{l} d[3,3] + \text{insert}(t[4]) = 3 \\ d[3,2] + \text{substitute}(s[3], t[4]) = 3 \\ d[4,2] + \text{delete}(s[3]) = 3 \end{array} \right\} = 3$$

最小编辑距离

stefanie

## 最小编辑距离的动态规划

- 原始串 `s o t`
- 目标串 `s t o p`

编辑操作①  
`s o t`  
 $\downarrow$   
`s t o t` (1. 插入t, 1分, 累计1分)  
 $\downarrow$   
`s t o p` (2. t替换p, 2分, 累计3分)

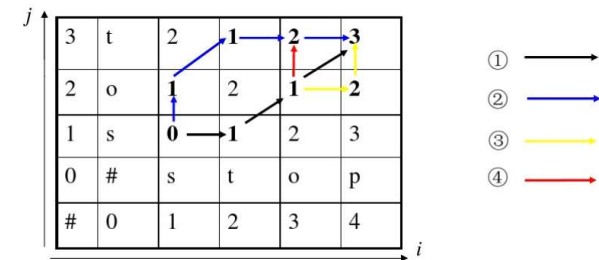
编辑操作②  
`s o t`  
 $\downarrow$   
`s t t` (1. 删除o, 1分, 累计1分)  
 $\downarrow$   
`s t o` (2. 插入o, 1分, 累计2分)  
 $\downarrow$   
`s t o p` (3. 插入p, 1分, 累计3分)

编辑操作③  
`s o t`  
 $\downarrow$   
`s t o t` (1. 插入t, 1分, 累计1分)  
 $\downarrow$   
`s t o p t` (2. 插入p, 1分, 累计2分)  
 $\downarrow$   
`s t o p` (3. 删除t, 1分, 累计3分)

编辑操作④  
`s o t`  
 $\downarrow$   
`s t o t` (1. 插入t, 1分, 累计1分)  
 $\downarrow$   
`s t o` (2. 删除t, 1分, 累计2分)  
 $\downarrow$   
`s t o p` (3. 插入p, 1分, 累计3分)

## 最小编辑距离的动态规划

- 原始串 `s o t`
- 目标串 `s t o p`



stefanie

## 最小编辑距离的动态规划

### 练习:

intention  $\longrightarrow$  execution

stefanie

## 最小编辑距离的动态规划

### 练习:

intention  $\longrightarrow$  execution

i n t e n t i o n  
↓ ↓ ↓ ↓ ↓  
e x e c u t i o n  
s s s s s  
2 2 2 2 2 = 10

i n t e n \* t i o n  
↓ ↓ ↓ ↓ ↓  
\* e x e c u t i o n  
d s s s i  
1 2 2 2 1 = 8

stefanie

## 其他应用

### 生物信息学中的最小编辑距离（相似度）

- 问题：找到以下两个序列中的对齐序列，其可能是核苷酸或者蛋白质的结构？

```
AGGCTATCACCTGACCTCCAGGCCGATGCCC
TAGCTATCACGACCGCGGTCGATTGCCCCGAC
```

- 要求得到以下对齐序列：

```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---
TAG-CTATCAC--GACCGC--GGTCGATTGCCCCGAC
```

Thank You ! 