

# Регулярные выражения

# Регулярные языки Клини

Пусть  $\Sigma$  – алфавит, тогда

- Пустой язык  $\emptyset$  является регулярным
- Язык  $\{\varepsilon\}$ , состоящий из пустого слова  $\varepsilon$ , является регулярным
- Для всех  $a \in \Sigma$  язык, состоящий из одной буквы  $\{a\}$ , регулярный
- Если  $L$ ,  $L_1$  и  $L_2$  - регулярные языки, то регулярными будут

- объединение  $L_1 \cup L_2$

$$\Sigma = \{a, b\}$$

$$L_1 = \{\varepsilon, bb, ba\}$$

$$L_2 = \{a, aab\}$$

$$L_1 \cup L_2 = \{\varepsilon, bb, ba, a, aab\}$$

# Регулярные языки Клини

Пусть  $\Sigma$  – алфавит, тогда

- Пустой язык  $\emptyset$  является регулярным
- Язык  $\{\varepsilon\}$ , состоящий из пустого слова  $\varepsilon$ , является регулярным
- Для всех  $a \in \Sigma$  язык, состоящий из одной буквы  $\{a\}$ , регулярный
- Если  $L$ ,  $L_1$  и  $L_2$  - регулярные языки, то регулярными будут

- объединение  $L_1 \cup L_2$

- конкатенация  $L_1 \bullet L_2$

$$\Sigma = \{a, b\}$$

$$L_1 = \{\varepsilon, bb, ba\}$$

$$L_2 = \{a, aab\}$$

$$L_1 \bullet L_2 = \{a, aab, bba, bbaab, baa, baaab\}$$

# Регулярные языки Клини

Пусть  $\Sigma$  – алфавит, тогда

- Пустой язык  $\emptyset$  является регулярным
- Язык  $\{\varepsilon\}$ , состоящий из пустого слова  $\varepsilon$ , является регулярным
- Для всех  $a \in \Sigma$  язык, состоящий из одной буквы  $\{a\}$ , регулярный
- Если  $L$ ,  $L_1$  и  $L_2$  – регулярные языки, то регулярными будут

- объединение  $L_1 \cup L_2$

$$\Sigma = \{a, b\}$$

- конкатенация  $L_1 \bullet L_2$

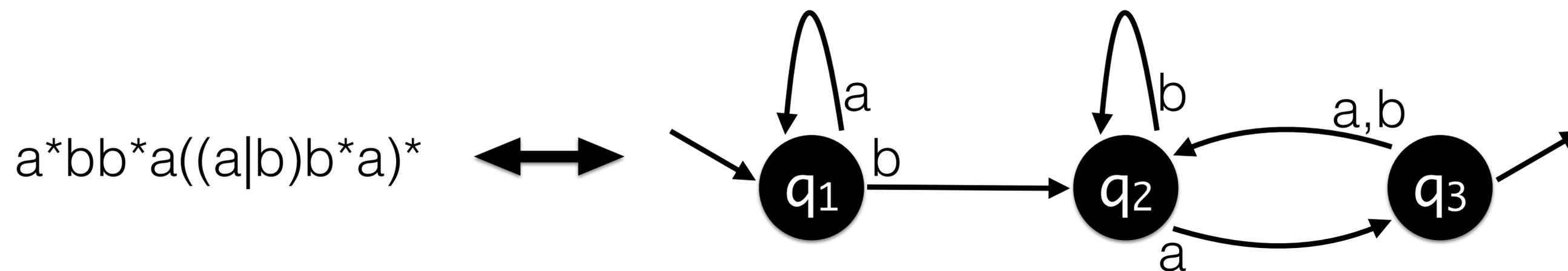
$$L = \{ab, b\}$$

$$L^* = \{\varepsilon, ab, b, abab, bab, abb, bb, (ab)^3, \dots\}$$

- результат применения звезды Клини  $L^*$

# Регулярные языки Клини

Регулярный язык  $\longleftrightarrow$  Детерминированный конечный автомат



# Регулярные выражения в Python

# Проблемы записи регулярных выражений

```
len('\\') == 1  
len(r'\\') == 2  
len('\\\\') == 2
```

Raw strings

# Задачи, решаемые регулярными выражениями

- Проверка соответствия строки шаблону
- Извлечение данных из строки по шаблону
- Изменение данных в строке, подходящих под шаблон

Модуль **re**

# Специальные символы в записи регулярных выражений

. \* + ? { } [ ] | ( ) ^ \$ \

# Специальные символы в записи регулярных выражений

- **.**  $\cdot \sim a, b, c, \dots \setminus \backslash n$  (re.DOTALL)
- **\***  $a^* \sim \varepsilon, a, aa, aaa, \dots; ab^* \sim a, ab, abb, \dots; (ab)^* \sim \varepsilon, ab, abab$
- **+**  $a+ \equiv aa^*$
- **?**  $ab?c \sim ac, abc$

# Про жадность

- \* → \*?
- + → +?
- ? → ??

# Специальные символы в записи регулярных выражений

- $\{m\}$   $a\{3\} \sim aaa$
- $\{m, n\}$   $a\{2, 4\} \sim aa, aaa, aaaa$
- $\{m, \}$   $a\{2, \} \equiv aa^+$
- $\{, n\}$   $a\{, 3\} \sim \varepsilon, a, aa, aaa$
- $[...]$   $[abc] \sim a, b, c; [0-9] \sim 0, 1, \dots, 9; [_0-9a-zA-Z]; [^abc]$
- $|$   $abc | de | f \sim abc, de, f$

# Специальные символы в записи регулярных выражений

- `()`
- `^`
- `$`

# Специальные символы в записи регулярных выражений

$\backslash d \equiv [0-9]$ ,  $\backslash D \equiv [^0-9]$

$\backslash s \equiv [ \backslash t \backslash n \backslash r \backslash f \backslash v ]$ ,  $\backslash S \equiv [^ \backslash s]$

- $\backslash \dots$

$\backslash w \equiv [a-zA-Z0-9_]$ ,  $\backslash W \equiv [^ \backslash w]$

$\backslash 1$ ,  $\backslash 2$ ,  $\backslash 3$ , ...

$\backslash \backslash$ ,  $\backslash .$ ,  $\backslash *$ , ...

# Специальные символы в записи регулярных выражений

- **(?:...)** non-grouping
- **(?=...)** positive lookahead
- **(?!...)** negative lookahead
- **(?<=...)** positive lookbehind
- **(?<!...)** negative lookbehind

# Вопросы производительности

```
for string in very_big_data:  
    re.search(r'something interesting', string)  
re.search(r'something interesting', string)
```

```
regex = re.compile(r'something interesting')  
for string in very_big_data:  
    regex.search(string)  
regex.search(string)
```

Извлечение данных из строки  
по шаблону

Изменение данных в строке,  
подходящих под шаблон