

# 预处理

原始文本 → 清洗 → 标准化 → 分词

网页  
新闻  
报告

无用标签  
特殊符号  
停用词

词干提取  
词形还原  
否定词处理  
敏感词处理

# 数据收集

1. 语料库：存放语言材料的数据库

类型：专用 & 通用 单语 & 多语 共时 & 历时 生语 & 熟语

2. 词汇知识库 (WordNet)：一个按语义关系网络组织的巨大语库

以同义词集合 (synset) 作为基本构建单位，将名词、动词、形容词、副词都组织到 synset 中，一个 synset 只包含一个解释，对一个 synset 中不同的词，分别用适当的例句加以区分。

如：car auto automobile 都属于“汽车”这个 synset

## WordNet 中的语义关系及标记符号

| 名词              |   | 动词                 |   | 形容词                       |   | 副词                   |   |
|-----------------|---|--------------------|---|---------------------------|---|----------------------|---|
| 反义关系<br>Antonym | ! | 反义关系<br>Antonym    | ! | 反义关系<br>Antonym           | ! | 反义关系<br>Antonym      | ! |
| 下位关系<br>Hyponym | ~ | 下位关系<br>Troponym   | ~ | 近义关系<br>Similar           | & | 导出形式<br>Derived from | \ |
| 上位关系<br>Hypemym | @ | 上位关系<br>Hypemym    | @ | 关系型形容词<br>Relational Adj. | \ |                      |   |
| 部分关系<br>Meronym | # | 蕴涵关系<br>Entailment | * | 又见<br>Also See            | ^ |                      |   |
| 整体关系<br>Holonym | % | 致使关系<br>Cause      | > | 属性<br>Attribute           | = |                      |   |
| 属性<br>Attribute | = | 又见<br>Also See     | ^ |                           |   |                      |   |

上下位关系：上位：抽象包含性概念  
下位：具体概念

如：animal ↔ dog

整体部分关系：car ↔ wheel

3. 数据爬虫

一种自动提取网页的程序，传统爬虫从一个或若干初始网页的 URL 开始，获得初始网页上的 URL，在抓取网页的过程中，不断从当前页面上抽取新的 URL 放入队列，直到满足一定条件。

步骤：① 获取网页：根据URL获取HTML数据

② 解析网页：解析HTML，获取HTML目标信息

③ 存储数据

## 西文字符编码

ASCII码 (美国标准信息交换码)

ASCII可显示字符 (共95个)

| 二进制       | 十进制 | 十六进制 | 图形      | 二进制       | 十进制 | 十六进制 | 图形 | 二进制       | 十进制 | 十六进制 | 图形 |
|-----------|-----|------|---------|-----------|-----|------|----|-----------|-----|------|----|
| 0010 0000 | 32  | 20   | (space) | 0100 0000 | 64  | 40   | @  | 0110 0000 | 96  | 60   | `  |
| 0010 0001 | 33  | 21   | !       | 0100 0001 | 65  | 41   | A  | 0110 0001 | 97  | 61   | a  |
| 0010 0010 | 34  | 22   | "       | 0100 0010 | 66  | 42   | B  | 0110 0010 | 98  | 62   | b  |
| 0010 0011 | 35  | 23   | #       | 0100 0011 | 67  | 43   | C  | 0110 0011 | 99  | 63   | c  |
| 0010 0100 | 36  | 24   | \$      | 0100 0100 | 68  | 44   | D  | 0110 0100 | 100 | 64   | d  |
| 0010 0101 | 37  | 25   | %       | 0100 0101 | 69  | 45   | E  | 0110 0101 | 101 | 65   | e  |
| 0010 0110 | 38  | 26   | &       | 0100 0110 | 70  | 46   | F  | 0110 0110 | 102 | 66   | f  |
| 0010 0111 | 39  | 27   | '       | 0100 0111 | 71  | 47   | G  | 0110 0111 | 103 | 67   | g  |
| 0010 1000 | 40  | 28   | (       | 0100 1000 | 72  | 48   | H  | 0110 1000 | 104 | 68   | h  |
| 0010 1001 | 41  | 29   | )       | 0100 1001 | 73  | 49   | I  | 0110 1001 | 105 | 69   | i  |
| 0010 1010 | 42  | 2A   | *       | 0100 1010 | 74  | 4A   | J  | 0110 1010 | 106 | 6A   | j  |
| 0010 1011 | 43  | 2B   | +       | 0100 1011 | 75  | 4B   | K  | 0110 1011 | 107 | 6B   | k  |
| 0010 1100 | 44  | 2C   | ,       | 0100 1100 | 76  | 4C   | L  | 0110 1100 | 108 | 6C   | l  |
| 0010 1101 | 45  | 2D   | -       | 0100 1101 | 77  | 4D   | M  | 0110 1101 | 109 | 6D   | m  |
| 0010 1110 | 46  | 2E   | .       | 0100 1110 | 78  | 4E   | N  | 0110 1110 | 110 | 6E   | n  |
| 0010 1111 | 47  | 2F   | /       | 0100 1111 | 79  | 4F   | O  | 0110 1111 | 111 | 6F   | o  |
| 0011 0000 | 48  | 30   | 0       | 0101 0000 | 80  | 50   | P  | 0111 0000 | 112 | 70   | p  |
| 0011 0001 | 49  | 31   | 1       | 0101 0001 | 81  | 51   | Q  | 0111 0001 | 113 | 71   | q  |
| 0011 0010 | 50  | 32   | 2       | 0101 0010 | 82  | 52   | R  | 0111 0010 | 114 | 72   | r  |
| 0011 0011 | 51  | 33   | 3       | 0101 0011 | 83  | 53   | S  | 0111 0011 | 115 | 73   | s  |
| 0011 0100 | 52  | 34   | 4       | 0101 0100 | 84  | 54   | T  | 0111 0100 | 116 | 74   | t  |
| 0011 0101 | 53  | 35   | 5       | 0101 0101 | 85  | 55   | U  | 0111 0101 | 117 | 75   | u  |
| 0011 0110 | 54  | 36   | 6       | 0101 0110 | 86  | 56   | V  | 0111 0110 | 118 | 76   | v  |
| 0011 0111 | 55  | 37   | 7       | 0101 0111 | 87  | 57   | W  | 0111 0111 | 119 | 77   | w  |
| 0011 1000 | 56  | 38   | 8       | 0101 1000 | 88  | 58   | X  | 0111 1000 | 120 | 78   | x  |
| 0011 1001 | 57  | 39   | 9       | 0101 1001 | 89  | 59   | Y  | 0111 1001 | 121 | 79   | y  |
| 0011 1010 | 58  | 3A   | :       | 0101 1010 | 90  | 5A   | Z  | 0111 1010 | 122 | 7A   | z  |
| 0011 1011 | 59  | 3B   | ;       | 0101 1011 | 91  | 5B   | [  | 0111 1011 | 123 | 7B   | {  |
| 0011 1100 | 60  | 3C   | <       | 0101 1100 | 92  | 5C   | \  | 0111 1100 | 124 | 7C   |    |
| 0011 1101 | 61  | 3D   | =       | 0101 1101 | 93  | 5D   | ]  | 0111 1101 | 125 | 7D   | }  |
| 0011 1110 | 62  | 3E   | >       | 0101 1110 | 94  | 5E   | ^  | 0111 1110 | 126 | 7E   | ~  |
| 0011 1111 | 63  | 3F   | ?       | 0101 1111 | 95  | 5F   | _  |           |     |      |    |

@hugachao123

## 中文字符编码

GB2312-80: 1980年发行的国标码, 覆盖常用简体字

GBK: 1995年发行, 覆盖生僻字、繁体字, windows默认中文编码, 国标扩展码

GB18030: 覆盖少数民族文字, 中国官方强制标准

GB13000/Unicode/ISO 10646: 国际通用字符集标准

GB2312-80: 用两个字节表示一个汉字, 每个字节的ASCII码都在161(A1) ~ 254(FF)间, 共 $94^2 = 8836$ 个码位, 共7445个字符(汉字6763个), 1391个空位

| 首字节<br>ASCII码 |                           | 161 | 254   |
|---------------|---------------------------|-----|-------|
| 161           | 标点、一般符号 (202)             | 161 | 非汉字区  |
| 162           | 序号 (60)                   | 168 |       |
| 163           | 数字 (22)、拉丁字母 (52)         |     |       |
| 164           | 日文假名 (169)                | 176 | 一级汉字区 |
| 165           |                           |     |       |
| 166           | 希腊字母 (48)                 |     |       |
| 167           | 俄文字母 (66)                 | 215 |       |
| 168           | 汉语拼音符号 (26)               | 216 | 二级汉字区 |
|               | 汉语注音字母 (37)               |     |       |
| 176~215       | 一级汉字 (3755个)<br>按汉语拼音顺序排序 |     |       |
| 216~247       | 二级汉字 (3008个)<br>按部首排列     | 247 |       |
|               |                           | 254 |       |

一个字在方阵中的坐标称为其区位码。

一级汉字：按拼音字母排序，同音字按笔形  
顺序排序（横直撇点折）

二级汉字：按部首顺序排序  
同部首按笔划数排序  
同笔划数按笔形顺序排序

Unicode：可视作 ISO/IEC 10646 标准的实践版

四维编码空间：

- 一个完整的 UCS-4 编码 = 4 个 8 位字节，记作：  
 $B_1$  (最高字节)  $B_2$  (次高字节)  $B_3$  (第3字节)  $B_4$  (最低字节)
- 根据最高位为 0 的最高字节分为  $2^1 = 128$  组 (group)  
即： $B_1$  0xxxxxxx
- 每个组再根据次高字节  $B_2$  分为  $2^8 = 256$  个平面 (plane)
- 每个平面根据第3字节  $B_3$  分为  $2^8 = 256$  行 (row)
- 每行有  $2^8$  个码位 (cell)

group 0 的 plane 0 被称作 BMP，是目前最实用的  
Unicode 版本，剩下 2 字节 ( $2^{16}$ ) 用于给字符编码。

Unicode 的实现方式不同于编码方式，因为一个字符  
的 Unicode 编码为 4 字节，会浪费大量空间。

Unicode的实现方式称为Unicode转换格式(UTF)  
UTF-8: 以无符号8位整数作编码单位  
类似的有UTF-16、UTF-32(32就是Unicode  
编码)

GBK: 兼容GB2312与Unicode

第一字节ASCII码: 129~254

第二字节ASCII码: 64~254(127除外)

## 字符编辑距离

定义: 在两个单词( $w_1, w_2$ )间由其中一个单词 $w_1$ 转换为另一个单词 $w_2$   
所需的最少单字符编辑操作次数(插入, 删除, 替换)

英文字符串编辑距离: 动态规划

字符串 $a[0 \dots m]$ 与 $b[0 \dots n]$ , 则 $d_{ij}$ 表示将 $a[0 \dots i-1]$ 转换为 $b[0 \dots j-1]$   
的最短编辑距离

$$\text{递推公式: } d_{ij} = \begin{cases} d_{i-1, j-1} & a_{i-1} = b_{j-1} \\ \min \begin{cases} d_{i-1, j} + \text{cost\_del} & \text{删 } a_{i-1} \\ d_{i, j-1} + \text{cost\_ins} & \text{插 } b_{j-1} \\ d_{i-1, j-1} + \text{cost\_rep} & \text{替} \end{cases} & a_{i-1} \neq b_{j-1} \end{cases}$$

汉字编辑距离: 由汉字笔顺的数字编码得来, 如“狼”与“狠”

“狼”: 3534511534    “狠”: 353511534

即该数字串编辑距离

## 正则表达式

通常用于检索、替换符合某个模式(规则)的文本

1. 普通字符: 直接匹配自身, 如“123”匹配“123”

2. 元字符

· : 匹配任意单个字符 如 "a.b" 匹配 "axb" "a#b" "a|b"

^ : 匹配字符串开头 (多行模式下匹配每行开头)

如: ^hello 匹配 "hello world" 但不匹配 "lhello"

\$ : 匹配字符串结尾 (或每行结尾)

如: world\$ 匹配 "hello world"

\* : 匹配前面的元素 0 或很多次

如: ab\* 匹配 a, ab, abb ...

+ : 匹配前面的元素 1 或很多次

如: ab+ 匹配 ab, abb, ...

? : 匹配前面的元素 0 或 1 次

如: ab? 匹配 a 和 ab

{n} : 匹配前面的元素 n 次 (n 为非负整数)

如: a{3} 匹配 aaa

{n,} : 匹配前面的元素至少 n 次

如: a{2,} 匹配 aa, aaa, ...

{n,m} : 匹配前面的元素 n 到 m 次

如: a{2,3} 匹配 aa, aaa

| : "或", 优先级最低的元素, 匹配左边或右边的模式

如: ab|cd 匹配 ab 或 cd

a(b|c)d 匹配 abd 或 acd

() : 分组, 将多个字符看作一个整体, 同时捕获匹配结果 (可以通过反向引用使用)

分组: (ab)+ 匹配 ab, abab ...

捕获与引用: (1d)\1 1d 匹配一个数字, 捕获为 \1  
\1 再使用一次刚捕获的数字  
匹配 11, 22, 33 ...

\: 转义字符, 取消元字符的特殊含义, 或表示预定义字符

如:  $a \backslash . b$  匹配  $a.b$

[ ]: 字符类, 匹配括号内任意字符, 支持范围表示 ( $[A-Z]$ ,  $[0-9]$ ), 加上  $^$  表否定

如:  $[abc]$  匹配  $a, b, c$

$[^abc]$  匹配任意非  $abc$  的字符

[ - ]: 字符类内的范围符, 表示连续字符集

如  $[a-z][0-9][a-zA-Z0-9]$

### 3. 预定义字符类

\d: 等价于  $[0-9]$

\D: 等价于  $[^0-9]$

\w: 等价于  $[a-zA-Z0-9_]$

注意这里有个下划线

\W: 等价于  $[^a-zA-Z0-9_]$

\s: 等价于  $[\t \n \r \f \v]$  匹配任意空白字符

\S: 等价于  $[^ \t \n \r \f \v]$

\b: 匹配单词边界 (单词与非单词分隔处, 如  $hello world$  中间空格两侧)

如:  $\backslash b \text{cat} \backslash b$  匹配  $\text{cat}$

\B: 匹配非单词边界

如:  $\text{cat} \backslash B$  不匹配  $\text{cat}$ , 但会匹配  $\text{catch}$  等

\n: 匹配换行符 (部分语言支持匹配  $\backslash r \backslash n$ )

\t: 匹配制表符

### 核心特性

#### 1. 贪婪匹配

默认情况下, 正则的量词 ( $* + ? \{n, m\}$ ) 是贪婪的, 即总是尽可能匹配更多字符, 在量词后面加  $?$  则反之

如: 匹配字符串 "aabbabcc"

a.\*b: 匹配 aabbab

a.\*?b: 匹配 aab

a.+c: 匹配 aabbabcc

a.+?c: 匹配 aabbabc

## 2. 分组与捕获

前面说 ( ) 会捕获匹配到的内容, 捕获的内容按捕获顺序 (\1, \2, ...) 编号, 可以被反向引用

如: (\w+)\s+\1: ( ) 会捕获 \w+ 匹配的字符串 (这里是一个单词) 并存入 \1, \s+ 是空白字符, \1 即引用 \1 里存的字符串, 如: "hello hello"

(\d{4}) - (\d{2}) - (\d{2}): 匹配 2026-04-07  
\1 存 "2026" \2 存 "04" \3 存 "07"

捕获自然会提升机器负担, 若不要捕获但要分组, 可以用: (?: ) , 如: (?: ab)

如果要为捕获组命名, 则用 (?<name> ) , 但只有部分语言支持 (Py, JS, Java)

## 3. 零宽断言

① 先行断言: 判断当前位置后面是否满足条件

· 前向先行: (?= pattern) 后面必须是 pattern

· 后向先行: (?< pattern) 后面必须不是 pattern

如: [a-z](?=.\*\d) 匹配后面跟着数字的字母 "a1" 的 "a"

^(?=.\*\d).\{6,12}\$: 这是密码的正则表达式, 要求这个密码长度在 6~12 个字符之间, 必须包含一个数字。"^" 强制匹配从字符串开头开始; (?=.\*\d) 表示至少包含一个数字的字符串, .\{6,12} 规定匹配长度在 6~12 之间, \$ 标记字符串结尾

注: `*\d` 若直接作为正则表达式, 那它会匹配“最长的以数字结尾的字符串”, 如“abc”里的“ab2”但是作为前向断言里的 pattern, 它就稍微不一样, 它的意义就变成了只要 `*\d` 能去匹配这个 pattern 就说明条件成立, 如 `pattern = ab3cd`, `*\d` 能匹配 `ab3`, 则条件成立。换句话说只要 pattern 含数字就行

② 后行断言: 判断当前位置前面是否满足条件

· 正向后行: `(?<=pattern)` 前面必须是 pattern

· 负向后行: `(?<!pattern)` 前面必须不是 pattern

如: `(?<=\d)[a-z]` 前面是数字的字母 如“1a”的“a”

## 常用组合

· `*`: 匹配任意不换行字符串

`[s\S]`: 匹配任意字符串

`[^x]*`: 匹配到 x 为止, 不包含 x

`[^"']*`: 匹配双引号内的内容

`\b\w+\b`: 匹配任意完整单词

`\S+`: 匹配一串非空字符

| Pattern   | Matches                  |
|---|--------------------------|
| <code>0\d{2,3}-\d{7,8}</code>                               | 电话号码 0785-2557392        |
| <code>^1[3 4 5 7 8]\d{9}\$</code>                           | 手机号码 15602527361         |
| <code>(^\d{15}\$) (^(\d{17})(\d X x)\$)</code>              | 身份证号码 44063719971204258X |
| <code>^(\w+)\@([0-9a-zA-Z]+)(\.[a-zA-Z]{2,4}){1,2}\$</code> | 电子邮件 694431056@qq.com    |
| <code>^\[\u4e00-\u9fa5]{0,}\$</code>                        | 检验汉字 你我他                 |
| <code>\n\s*\r</code>  | 空白行                      |
| <code>[1-9][0-9]{4,}</code>                                 | QQ号 从10000开始             |

| Pattern           | Matches  |
|-------------------|--|
| 网址 (URL)          | <code>[a-zA-z]+://[^\s]*</code>  |
| IP地址 (IP Address) | <code>((2[0-4]\d 25[0-5] ([01]?\d\d?)\.)\{3\}(2[0-4]\d 25[0-5] ([01]?\d\d?))</code>  |
| 日期 (年-月-日)        | <code>(\d{4})\d{2})-((1[0-2]) (0?[1-9]))-(((12 [0-9]) (3[01]) (0?[1-9]))</code>  |
| 日期 (月/日/年)        | <code>((1[0-2]) (0?[1-9]))/(((12 [0-9]) (3[01]) (0?[1-9]))/(\d{4})\d{2})</code>  |
| 时间 (小时:分钟, 24小时制) | <code>((1 0?)\d{0-1})\d{0-1}:([0-5]\d{0-1})</code>   |
| 中文及全角标点符号 (字符)    | <code>[\u3000-\u301e\u2010-\u2019\u201c-\u201d\u201e-\u201f\u2020-\u2021\u2022-\u2023\u2024-\u2025\u2026-\u2027\u2028-\u2029\u202a-\u202b\u202c-\u202d\u202e-\u202f\u2030-\u2031\u2032-\u2033\u2034-\u2035\u2036-\u2037\u2038-\u2039\u203a-\u203b\u203c-\u203d\u203e-\u203f\u2040-\u2041\u2042-\u2043\u2044-\u2045\u2046-\u2047\u2048-\u2049\u204a-\u204b\u204c-\u204d\u204e-\u204f\u2050-\u2051\u2052-\u2053\u2054-\u2055\u2056-\u2057\u2058-\u2059\u205a-\u205b\u205c-\u205d\u205e-\u205f\u2060-\u2061\u2062-\u2063\u2064-\u2065\u2066-\u2067\u2068-\u2069\u206a-\u206b\u206c-\u206d\u206e-\u206f\u2070-\u2071\u2072-\u2073\u2074-\u2075\u2076-\u2077\u2078-\u2079\u207a-\u207b\u207c-\u207d\u207e-\u207f\u2080-\u2081\u2082-\u2083\u2084-\u2085\u2086-\u2087\u2088-\u2089\u208a-\u208b\u208c-\u208d\u208e-\u208f\u2090-\u2091\u2092-\u2093\u2094-\u2095\u2096-\u2097\u2098-\u2099\u209a-\u209b\u209c-\u209d\u209e-\u209f\u20a0-\u20a1\u20a2-\u20a3\u20a4-\u20a5\u20a6-\u20a7\u20a8-\u20a9\u20aa-\u20ab\u20ac-\u20ad\u20ae-\u20af\u20b0-\u20b1\u20b2-\u20b3\u20b4-\u20b5\u20b6-\u20b7\u20b8-\u20b9\u20ba-\u20bb\u20bc-\u20bd\u20be-\u20bf\u20c0-\u20c1\u20c2-\u20c3\u20c4-\u20c5\u20c6-\u20c7\u20c8-\u20c9\u20ca-\u20cb\u20cc-\u20cd\u20ce-\u20cf\u20d0-\u20d1\u20d2-\u20d3\u20d4-\u20d5\u20d6-\u20d7\u20d8-\u20d9\u20da-\u20db\u20dc-\u20dd\u20de-\u20df\u20e0-\u20e1\u20e2-\u20e3\u20e4-\u20e5\u20e6-\u20e7\u20e8-\u20e9\u20ea-\u20eb\u20ec-\u20ed\u20ee-\u20ef\u20f0-\u20f1\u20f2-\u20f3\u20f4-\u20f5\u20f6-\u20f7\u20f8-\u20f9\u20fa-\u20fb\u20fc-\u20fd\u20fe-\u20ff\u2010-\u2019\u201c-\u201d\u201e-\u201f\u2020-\u2021\u2022-\u2023\u2024-\u2025\u2026-\u2027\u2028-\u2029\u202a-\u202b\u202c-\u202d\u202e-\u202f\u2030-\u2031\u2032-\u2033\u2034-\u2035\u2036-\u2037\u2038-\u2039\u203a-\u203b\u203c-\u203d\u203e-\u203f\u2040-\u2041\u2042-\u2043\u2044-\u2045\u2046-\u2047\u2048-\u2049\u204a-\u204b\u204c-\u204d\u204e-\u204f\u2050-\u2051\u2052-\u2053\u2054-\u2055\u2056-\u2057\u2058-\u2059\u205a-\u205b\u205c-\u205d\u205e-\u205f\u2060-\u2061\u2062-\u2063\u2064-\u2065\u2066-\u2067\u2068-\u2069\u206a-\u206b\u206c-\u206d\u206e-\u206f\u2070-\u2071\u2072-\u2073\u2074-\u2075\u2076-\u2077\u2078-\u2079\u207a-\u207b\u207c-\u207d\u207e-\u207f\u2080-\u2081\u2082-\u2083\u2084-\u2085\u2086-\u2087\u2088-\u2089\u208a-\u208b\u208c-\u208d\u208e-\u208f\u2090-\u2091\u2092-\u2093\u2094-\u2095\u2096-\u2097\u2098-\u2099\u209a-\u209b\u209c-\u209d\u209e-\u209f\u20a0-\u20a1\u20a2-\u20a3\u20a4-\u20a5\u20a6-\u20a7\u20a8-\u20a9\u20aa-\u20ab\u20ac-\u20ad\u20ae-\u20af\u20b0-\u20b1\u20b2-\u20b3\u20b4-\u20b5\u20b6-\u20b7\u20b8-\u20b9\u20ba-\u20bb\u20bc-\u20bd\u20be-\u20bf\u20c0-\u20c1\u20c2-\u20c3\u20c4-\u20c5\u20c6-\u20c7\u20c8-\u20c9\u20ca-\u20cb\u20cc-\u20cd\u20ce-\u20cf\u20d0-\u20d1\u20d2-\u20d3\u20d4-\u20d5\u20d6-\u20d7\u20d8-\u20d9\u20da-\u20db\u20dc-\u20dd\u20de-\u20df\u20e0-\u20e1\u20e2-\u20e3\u20e4-\u20e5\u20e6-\u20e7\u20e8-\u20e9\u20ea-\u20eb\u20ec-\u20ed\u20ee-\u20ef\u20f0-\u20f1\u20f2-\u20f3\u20f4-\u20f5\u20f6-\u20f7\u20f8-\u20f9\u20fa-\u20fb\u20fc-\u20fd\u20fe-\u20ff\u2010-\u2019\u201c-\u201d\u201e-\u201f\u2020-\u2021\u2022-\u2023\u2024-\u2025\u2026-\u2027\u2028-\u2029\u202a-\u202b\u202c-\u202d\u202e-\u202f\u2030-\u2031\u2032-\u2033\u2034-\u2035\u2036-\u2037\u2038-\u2039\u203a-\u203b\u203c-\u203d\u203e-\u203f\u2040-\u2041\u2042-\u2043\u2044-\u2045\u2046-\u2047\u2048-\u2049\u204a-\u204b\u204c-\u204d\u204e-\u204f\u2050-\u2051\u2052-\u2053\u2054-\u2055\u2056-\u2057\u2058-\u2059\u205a-\u205b\u205c-\u205d\u205e-\u205f\u2060-\u2061\u2062-\u2063\u2064-\u2065\u2066-\u2067\u2068-\u2069\u206a-\u206b\u206c-\u206d\u206e-\u206f\u2070-\u2071\u2072-\u2073\u2074-\u2075\u2076-\u2077\u2078-\u2079\u207a-\u207b\u207c-\u207d\u207e-\u207f\u2080-\u2081\u2082-\u2083\u2084-\u2085\u2086-\u2087\u2088-\u2089\u208a-\u208b\u208c-\u208d\u208e-\u208f\u2090-\u2091\u2092-\u2093\u2094-\u2095\u2096-\u2097\u2098-\u2099\u209a-\u209b\u209c-\u209d\u209e-\u209f\u20a0-\u20a1\u20a2-\u20a3\u20a4-\u20a5\u20a6-\u20a7\u20a8-\u20a9\u20aa-\u20ab\u20ac-\u20ad\u20ae-\u20af\u20b0-\u20b1\u20b2-\u20b3\u20b4-\u20b5\u20b6-\u20b7\u20b8-\u20b9\u20ba-\u20bb\u20bc-\u20bd\u20be-\u20bf\u20c0-\u20c1\u20c2-\u20c3\u20c4-\u20c5\u20c6-\u20c7\u20c8-\u20c9\u20ca-\u20cb\u20cc-\u20cd\u20ce-\u20cf\u20d0-\u20d1\u20d2-\u20d3\u20d4-\u20d5\u20d6-\u20d7\u20d8-\u20d9\u20da-\u20db\u20dc-\u20dd\u20de-\u20df\u20e0-\u20e1\u20e2-\u20e3\u20e4-\u20e5\u20e6-\u20e7\u20e8-\u20e9\u20ea-\u20eb\u20ec-\u20ed\u20ee-\u20ef\u20f0-\u20f1\u20f2-\u20f3\u20f4-\u20f5\u20f6-\u20f7\u20f8-\u20f9\u20fa-\u20fb\u20fc-\u20fd\u20fe-\u20ff</code> |
| 中国大陆邮政编码          | <code>[1-9]\d{5}</code>  |
| 整数                | <code>-?\d+</code>   |
| 小数                | <code>(-?\d+)(\.\d+)?</code>   |
| 不包含 abc 的单词       | <code>\b((?!abc)\w)+\b</code>  |