# Python快速入门

嵩天



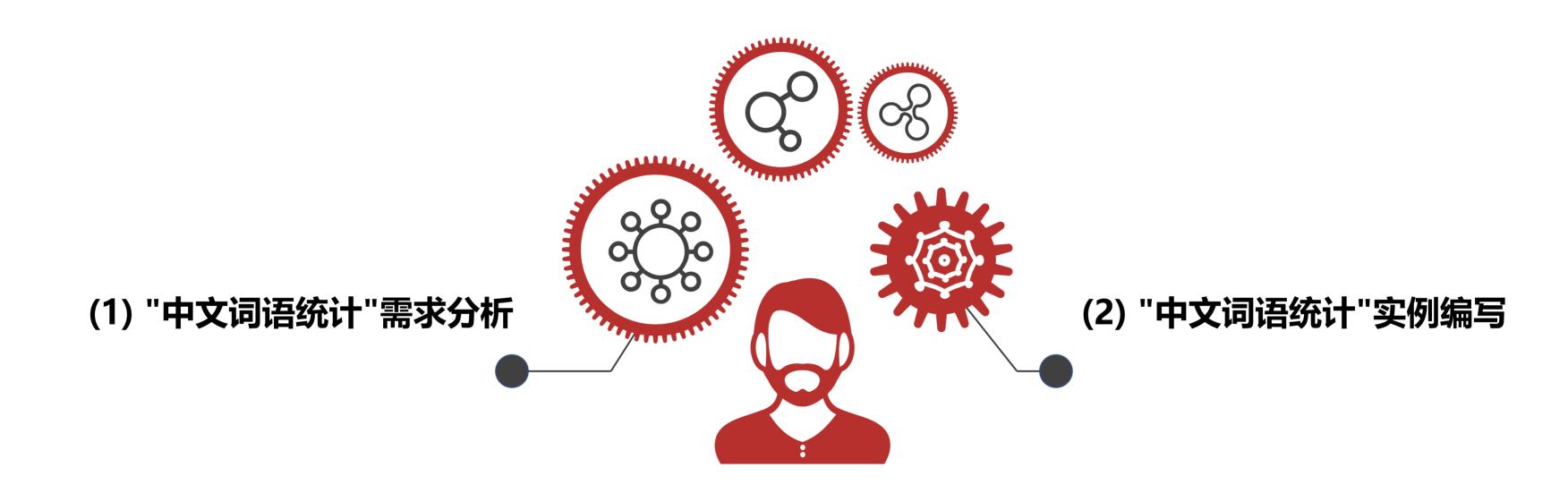
## 实例3:中文词语统计

嵩天



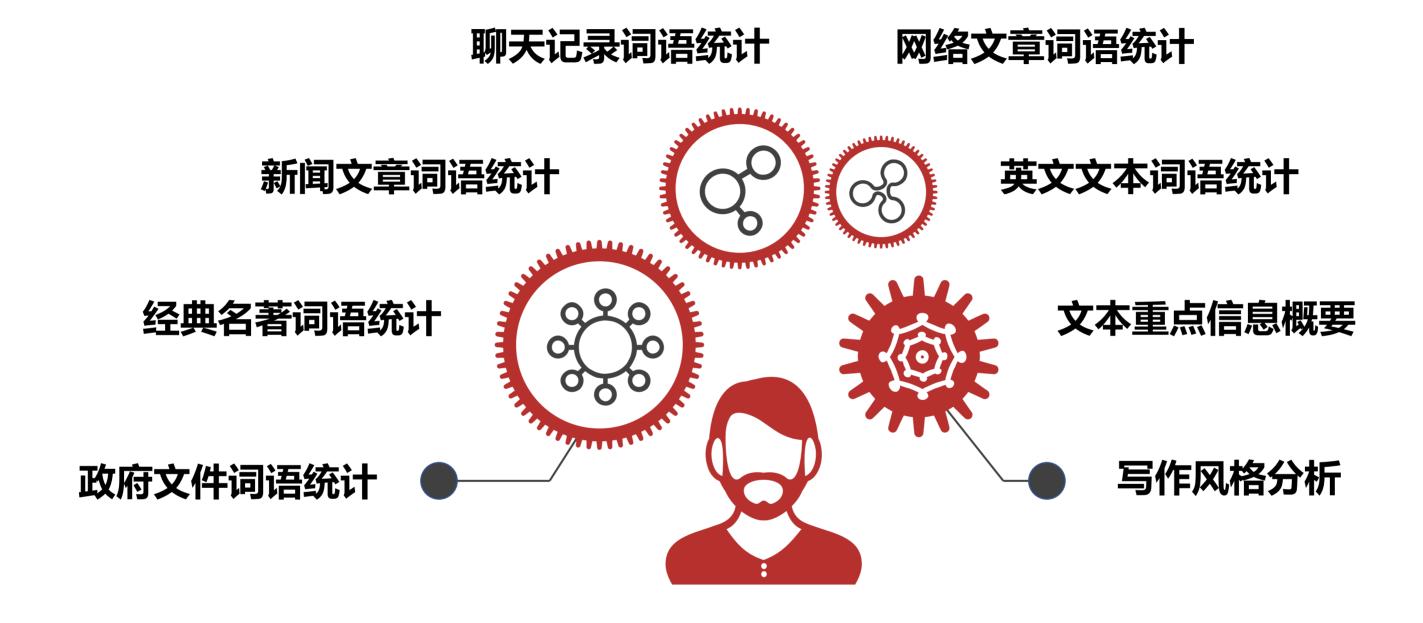


## 单元开篇



实例3:中文词语统计

## 单元开篇



实例3:中文词语统计



### 程序需求

#### 统计中文词语出现次数

- 以政府一号文件为例,统计出现的中文词语数量
- 按照一定标准输出,如出现次数等
- 需要解决中文分词问题,如:这是一门好课 -> 这是 一门 好课

## 程序需求

#### 统计中文词语出现次数

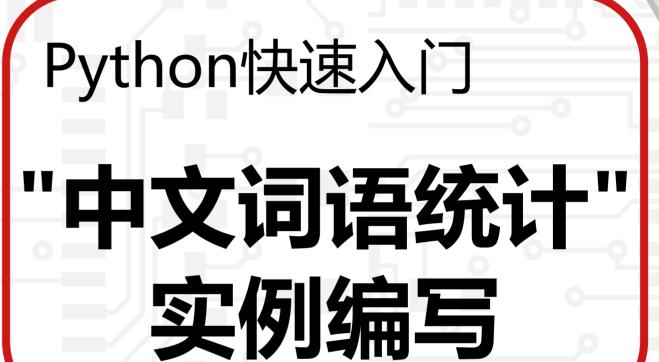
• 以每年政府一号文件为例

## 程序需求

#### 统计中文词语出现次数

• 输入: 2018年一号文件.txt

• 输出: 出现次数超过50次的词语, 不包括换行



```
#WordsCount.py
import jieba
f = open("2018年一号文件.txt", "r", encoding="utf-8")
txt = f.read()
f.close()
Is = jieba. lcut(txt)
d = \{\}
for w in s:
    d[w] = d. get(w, 0) + 1
for k in d:
    if d[k] >= 50 and k != "\n":
        print('"{}"出现{}次'.format(k, d[k]))
```

```
#WordsCount.py
     import jieba
     |f = open("2018年一号文件.txt", "r", encoding="utf-8")
     txt = f.read()
     f.close()
     Is = jieba. lcut(txt)
     d = \{\}
     for w in s:
         d[w] = d. get(w, 0) + 1
     for k in d:
         if d[k] >= 50  and k != "\n":
             print('"{}"出现{}次'.format(k, d[k]))
```

主释

```
#WordsCount.py
--> import jieba
    f = open("2018年一号文件.txt", "r", encoding="utf-8")
    txt = f.read()
    f.close()
     Is = jieba.lcut(txt)
    d = \{\}
     for w in s:
        d[w] = d. get(w, 0) + 1
     for k in d:
         if d[k] >= 50 and k != "\n":
            print('"{}"出现{}次'.format(k, d[k]))
```

引入外部 功能库

```
#WordsCount.py
    import jieba
打开文件
    txt = f.read()
f. close()
                                                  关闭文件
    Is = jieba. lcut(txt)
    d = \{\}
    for w in s:
       d[w] = d. get(w, 0) + 1
    for k in d:
       if d[k] >= 50 \ and \ k != "\n":
          print('"{}"出现{}次'.format(k, d[k]))
```



```
#WordsCount.py
      import jieba
     |f = open("2018年一号文件.txt", "r", encoding="utf-8")
\rightarrow txt = f. read()
     f.close()
      Is = jieba. lcut(txt)
     d = \{\}
     for w in s:
         d[w] = d. get(w, 0) + 1
      for k in d:
          if d[k] >= 50 and k != "\n":
              print('"{}"出现{}次'.format(k, d[k]))
```

读入文本

```
#WordsCount.py
     import jieba
    f = open("2018年一号文件.txt", "r", encoding="utf-8")
    txt = f.read()
    f.close()
Is = jieba. lcut(txt)
     d = \{\}
     for w in s:
         d[w] = d. get(w, 0) + 1
     for k in d:
         if d[k] >= 50  and k != "\n":
             print('"{}"出现{}次'.format(k, d[k]))
```

python

```
#WordsCount.py
     import jieba
     f = open("2018年一号文件.txt", "r", encoding="utf-8")
     txt = f.read()
     f.close()
     Is = jieba.lcut(txt)
\rightarrow d = \{\}
     for w in s:
         d[w] = d. get(w, 0) + 1
     for k in d:
         if d[k] >= 50  and k != "\n":
             print('"{}"出现{}次'.format(k, d[k]))
```

建立字典



```
#WordsCount.py
     import jieba
    f = open("2018年一号文件.txt", "r", encoding="utf-8")
    txt = f.read()
    f.close()
     Is = jieba. lcut(txt)
    d = \{\}
--- for w in Is:
d[w] = d. get(w, 0) + 1
     for k in d:
         if d[k] >= 50  and k != "\n":
            print('"{}"出现{}次'.format(k, d[k]))
```

利用字典 词语统计



```
#WordsCount.py
     import jieba
     f = open("2018年一号文件.txt", "r", encoding="utf-8")
     txt = f.read()
     f.close()
     Is = jieba. lcut(txt)
     d = \{\}
     for w in s:
         d[w] = d. get(w, 0) + 1
\rightarrow for k in d:
         if d[k] >= 50 and k != "\n":
             print('"{}"出现{}次'.format(k, d[k]))
```

遍历结果 设置条件 打印输出



## 注意事项

#### 相比其他编程语言

- 每行后没有分号;
- · 没有begin, end, {, }等表示代码归属的元素, 只用缩进表达代码所属关系
- 变量直接使用,无需类型声明
- import可以引入外部功能库



## Thank you