

数据清洗在网络安全中的应用

数据都不会清洗，你搭个毛线的社工库

演讲者：耿浩然

云南水熊科技有限公司

● 清洗数据的动机

1. 数据不正确或不一致会导致公共或者私人规模的错误结论和错误的行为
2. 减少人力成本
3. 提高数据应用效率

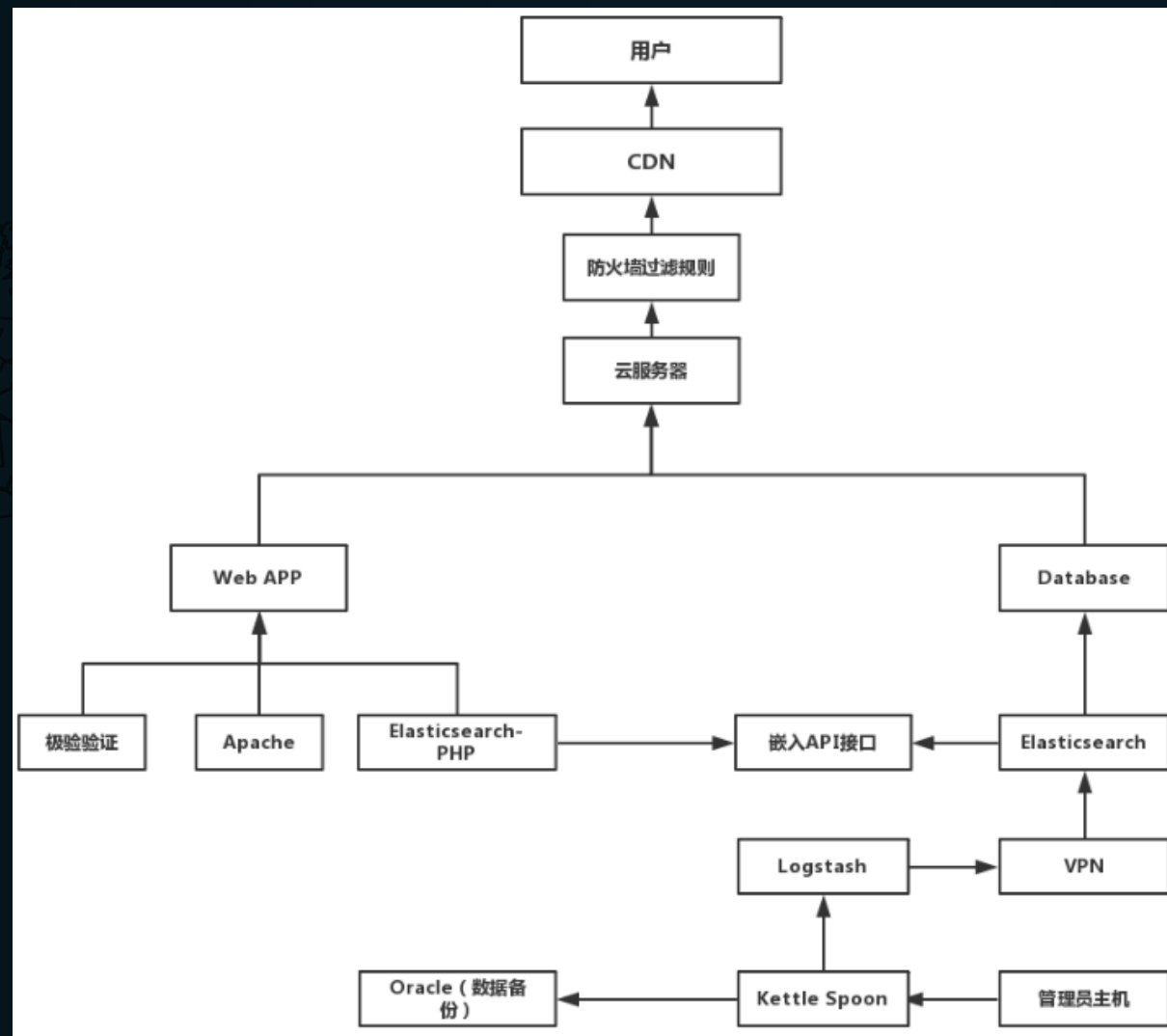


应用实例



● 某个社工库的技术实现图

ETL工具: Kettle Spoon
索引: Elasticsearch
入库: Logitech



● 数据分析

增：数据分析

删：删除无用字段

改：改变现有字段

查：要查询的字段

id	email	password	time
1	test@test.com	123	2017
2		3	2017
3	test@test.com		2017
3	test@test.com		2017
4	test@test.com		2017

● 模型构建

输入：起始数据全部转换为txt文本，采用统一分隔符，相同字段数

第一次过滤：过滤缺失数据

错误输出：将过滤的缺失数据输出到指定路径

第二次过滤：对数据MD5加密，然后进行排序去重

字段选择：删除和增加字段

输出：输出到logitech配置文件指定目录

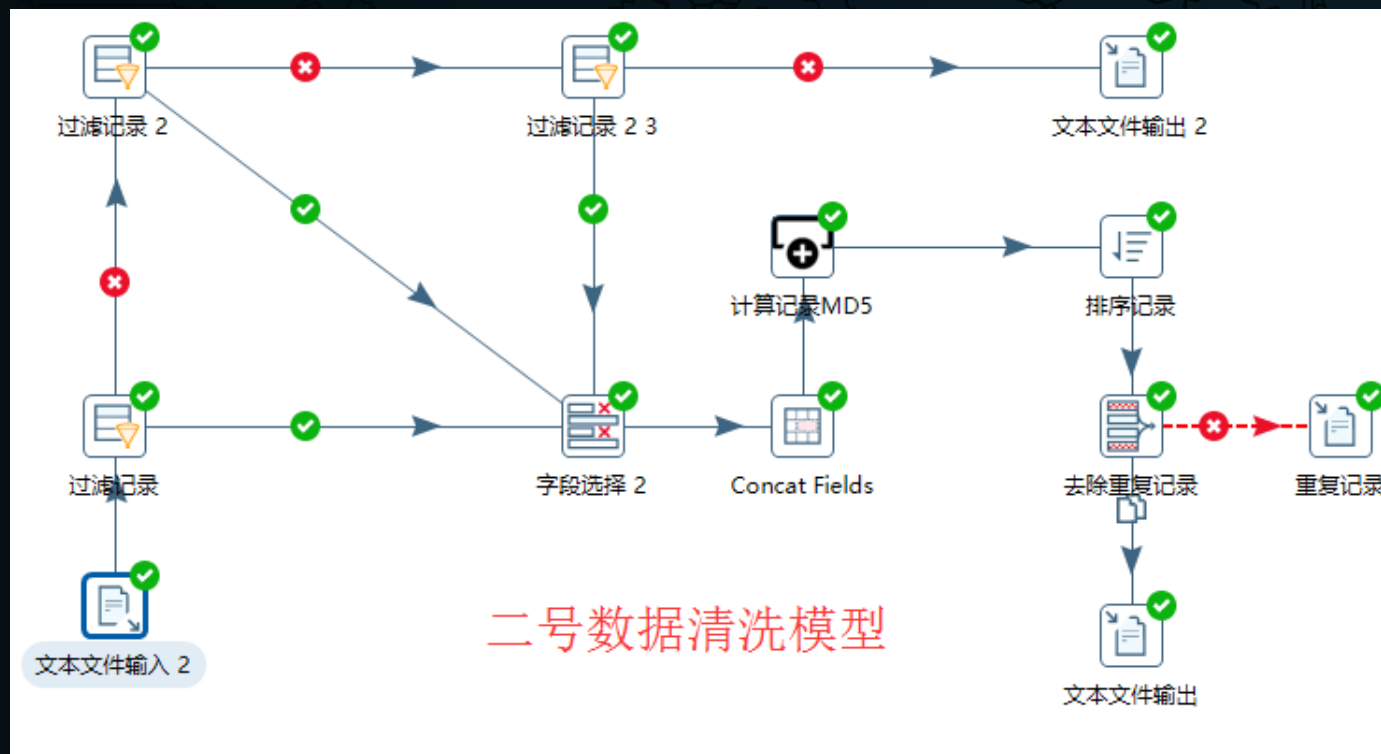
● Logitech config配置文件

```
1 input {
2   file {
3     path => ["E:/test/updata/*.txt"]
4     start_position => "beginning"
5     ignore_old => 999999999
6     codec => plain {charset => ["UTF-8"]}
7   }
8 }
9
10 filter {
11   csv {
12     separator => "|~|^|"
13     source => "message"
14     columns => ["username", "password", "email", "source"]
15   }
16 mutate { remove_field => ["message", "host", "@timestamp", "path", "@version"] }
17 }
18 output {
19   elasticsearch {
20     hosts => "127.0.0.1"
21     index => "test"
22     document_type => "test"
23     flush_size => "60000"
24   }
25   stdout {
26     codec => json
27   }
28 }
```

三秒钟模型演示之数据过滤

第一道过滤条件：条数据需要有username和password，username和email，或者password和email。

第二道过滤条件：对数据进行MD5加密排序后去重。



● 数据处理结构图概览

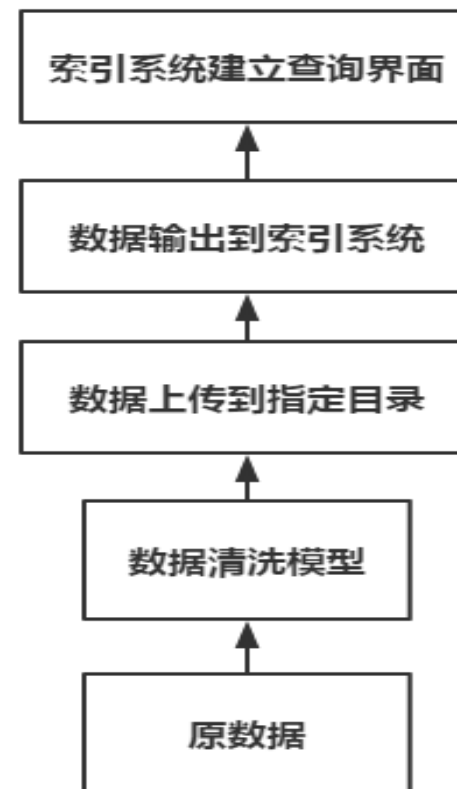
必要条件：

相同字段数量

字段value排序相同

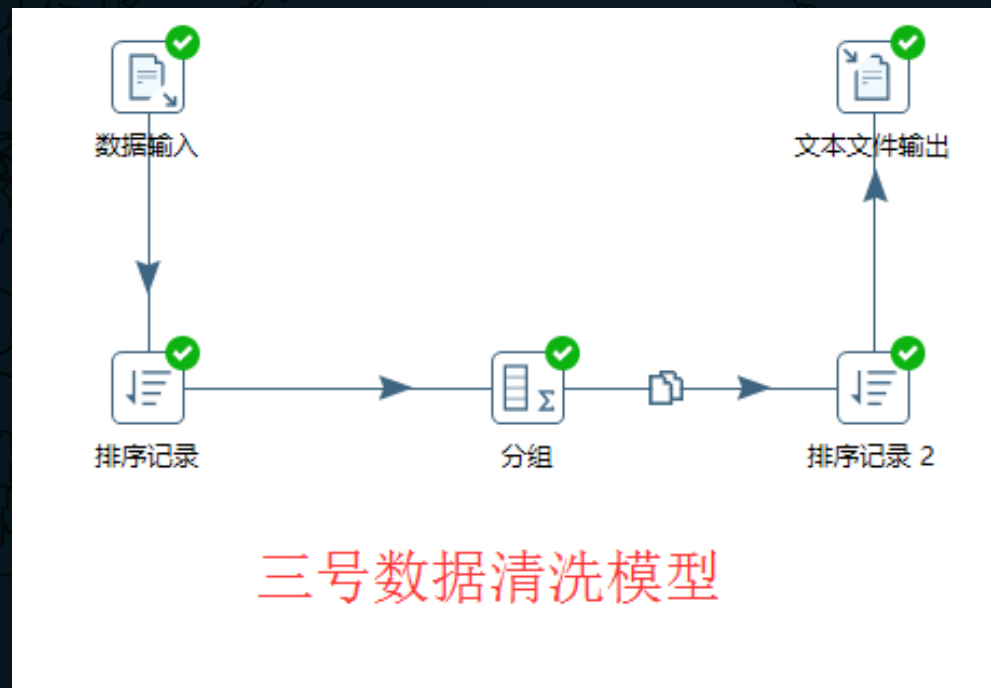
数据分隔符和封闭符相同

数据编码相同



一分钟模型演示之数据统计

弱口令字典清洗



● Reference list

<http://ieeexplore.ieee.org/document/7307098/> A Survey of Data Mining and Machine Learning Methods for Cyber Security Intrusion Detection

<http://www.aqniu.com/news-views/9369.html> BlackHat2015焦点访谈：中国黑客如何玩转深度学习

<http://forums.pentaho.com/showthread.php?59651-MD5-Hash> Thread: MD5 Hash

<http://wiki.pentaho.com/display/EAI/PDI+Spoon+Plugin+Development> PDI Spoon Plugin Development

谢谢观看

制作

By-耿浩然



18687117526



水熊科技-Arthur (耿浩然)

云南 昆明

