

# Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 6)

## Задание:

- 1) Установить виртуальную машину с ОС Ubuntu в VirtualBox.
- 2) Установить Elasticsearch, Neo4j, Hadoop, Spark.
- 3) Вручную создать JSON-файл с 20-30 JSON-документами для предметной области, указанной в варианте.
- 4) В Elasticsearch создать индекс с анализатором и маппингом, проиндексировать JSON-документы, разработать запросы с вложенной агрегацией, представить результаты в среде Kibana.
- 5) В Neo4j по данным из Elasticsearch заполнить графовую базу данных, разработать и реализовать запрос к этой БД.
- 6) В Spark по данным из Elasticsearch сформировать csv-файлы с таблицами и сохранить их в файловой системе HDFS, написать запрос и реализовать его в Spark, проанализировать процесс выполнения запроса с использованием монитора Spark.

## Предметная область – Салон красоты

### Elasticsearch:

1. Типы JSON-документов:

Процедура:

```
{index, doc_type, id, body: {сведения_о_процедуре*, стоимость}}
```

Пациент:

```
{index, doc_type, id, body: {id_пациента, персональные_данные*, путёвка, дата_прибытия, продолжительность_прибывания, [диагноз*], [id_процедуры]}}
```

2. Требования к анализатору: для полей, отмеченных символом \*, должен быть разработан анализатор со следующими требованиями: разделить текст на слова, убрать пунктуацию с помощью токенизатора standard (русский), перевести все токены в нижний регистр, убрать токены, находящиеся в списке стоп-слов, выполнить стемминг оставшихся токенов с помощью фильтра snowball.

3. Запросы с вложенной агрегацией:

- разбить пациентов по дате прибытия с периодом 1 год, для каждой «корзины» определить количество пациентов, прошедших каждую процедуру;
- определить стоимость процедуры по заданным ключевым словам.

### Neo4j:

1. По данным из Elasticsearch заполнить графовую базу данных Пациент(id\_пациента, дата\_прибытия, персональные\_данные) -Прошёл(стоимость) - Процедура(id\_процедуры).
2. Разработать и реализовать запрос: найти процедуру с максимальной суммарной стоимостью.

### Spark:

1. По данным из Elasticsearch сформировать csv-файлы (с внутренней схемой) таблиц «Пациент», «Назначение», «Процедура» и сохранить их в файловой системе HDFS.
2. Написать запрос select: определить суммарное число назначений по каждой процедуре.
3. Реализовать этот запрос в Spark. Построить временную диаграмму его выполнения по результатам работы монитора.

					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)							
					Задание на курсовой проект	Литер.			Масса		Масштаб	
Изм	Лист	№ документа	Подпись	Дата								
Разраб.												
Руковод.												
Н. Контр.						Лист 1			Листов 10			

# Индексация документов Elasticsearch

## Анализатор для индексов

```
"analysis" : {
  "filter": {
    "russian_stop_words": {
      "type": "stop",
      "stopwords": "_russian_"
    },
    "filter_ru_sn": {
      "type": "snowball",
      "language": "Russian"
    }
  },
  "analyzer": {
    "analitic_for_ru": {
      "type": "custom",
      "tokenizer": "standard",
      "filter": [
        "lowercase",
        "russian_stop_words",
        "filter_ru_sn"
      ]
    }
  }
}
```

## Маппинг для индекса процедура

```
ProcedureMapping = {
  "properties": {
    "name": {
      "type": "text",
      "fielddata": True,
      "analyzer": "analitic_for_ru",
      "search_analyzer": "analitic_for_ru"
    },
    "description": {
      "type": "text",
      "fielddata": True,
      "analyzer": "analitic_for_ru",
      "search_analyzer": "analitic_for_ru"
    },
    "price": {
      "type": "integer"
    }
  }
}
```

## Фрагмент маппинга для индекса пациент

```
PatientMapping = {
  "properties": {
    "patient_id": {
      "type": "text",
      "fielddata": True
    },
    "personal_data": {
      "type": "text",
      "analyzer": "analitic_for_ru",
      "search_analyzer": "analitic_for_ru",
      "fielddata": True
    },
    "precedures_id": {
      "type": "text",
      "fielddata": True
    }
  }
}
```

Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)

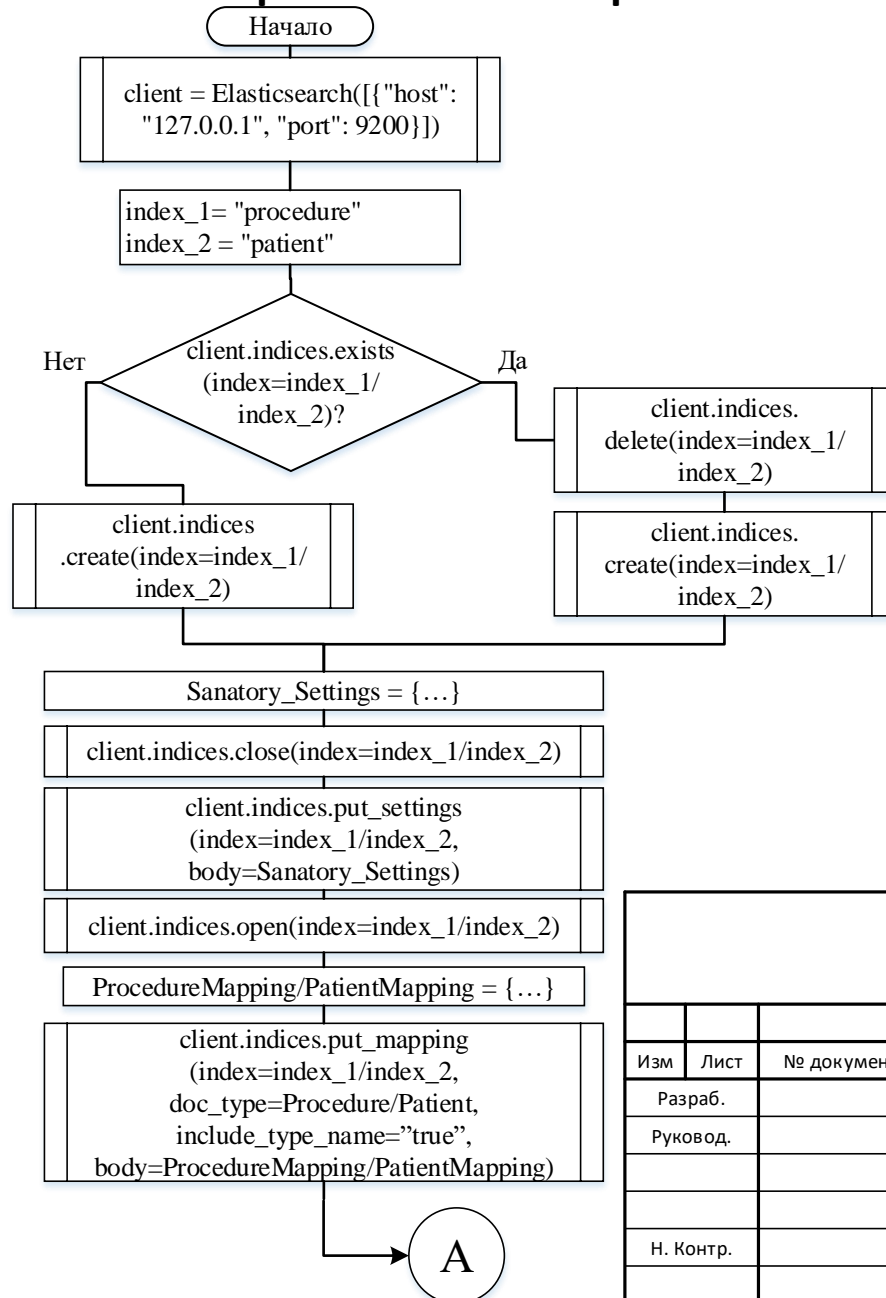
Изм.	Лист	№ документа	Подпись	Дата
Разраб.				
Руковод.				
Н. Контр.				

Индексация документов  
Elasticsearch. Маппинг и  
анализатор

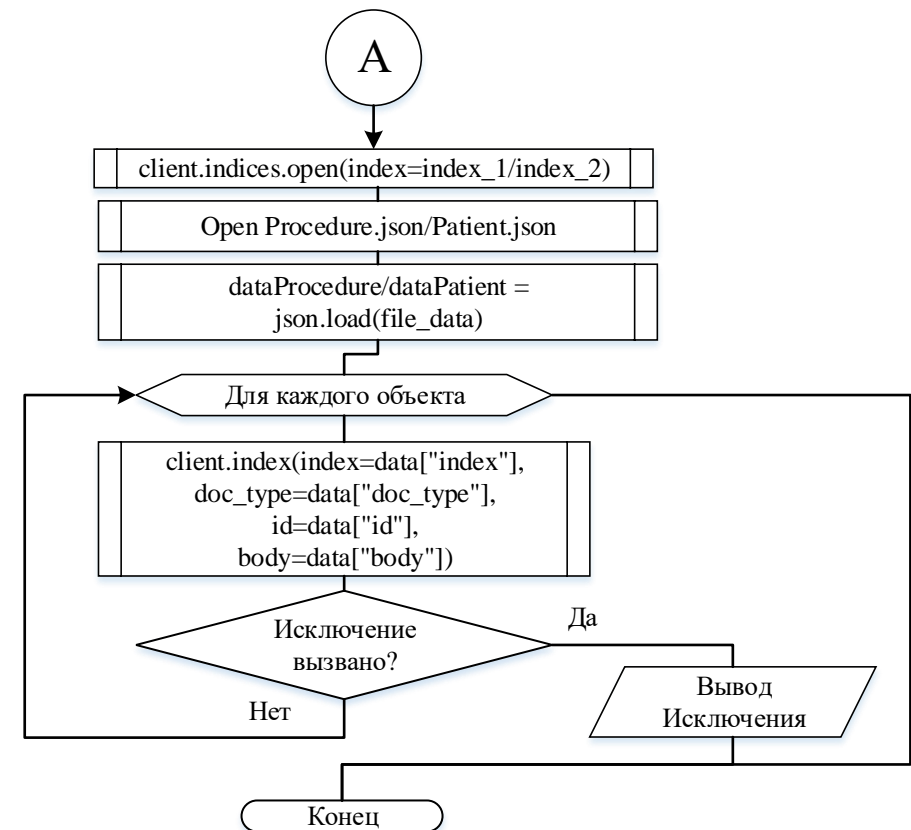
Литер.	Масса	Масштаб
Лист 2		Листов 10

# Индексация документов Elasticsearch

## Алгоритм добавления маппинга и настройки анализатора



## Алгоритм индексации



Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)

Индексация документов Elasticsearch. Алгоритмы индексации

Литер.			Масса	Масштаб
Лист 3			Листов 10	

# Elasticsearch. Запросы

**Первый запрос:** разбить пациентов по дате прибытия с периодом 1 год, для каждой «корзины» определить количество пациентов, прошедших каждую процедуру.

```
1- {
2  "took" : 27,
3  "timed_out" : false,
4  "_shards" : {
5    "total" : 1,
6    "successful" : 1,
7    "skipped" : 0,
8    "failed" : 0
9  },
10 "hits" : {
11   "total" : {
12     "value" : 30,
13     "relation" : "eq"
14   },
15   "max_score" : null,
16   "hits" : [ ]
17 },
18 "aggregations" : {
19   "year_period" : {
20     "buckets" : [
21       {
22         "key_as_string" : "2018-01-01",
23         "key" : 1514764800000,
24         "doc_count" : 8,
25         "Procedure" : {
26           "doc_count_error_upper_bound" : 0,
27           "sum_other_doc_count" : 17,
28           "buckets" : [
29             {
30               "key" : "pr01",
31               "doc_count" : 3,
32               "number_of_patients" : {
33                 "value" : 3
34               }
35             },
36             {
37               "key" : "pr02",
38               "doc_count" : 1,
39               "number_of_patients" : {
```

```
1 GET patient/_search
2 {
3   "size": 0,
4   "aggregations": {
5     "year_period": {
6       "date_histogram": {
7         "field": "date_of_arrival",
8         "calendar_interval": "year",
9         "format": "yyyy-MM-dd"
10      },
11      "aggregations": {
12        "Procedure": {
13          "terms": {
14            "field": "procedures_id",
15            "order": {
16              "key": "asc"
17            }
18          },
19          "aggregations": {
20            "number_of_patients": {
21              "value_count": {
22                "field": "patient_id"
23              }
24            }
25          }
26        }
27      }
28    }
29  }
30 }
```

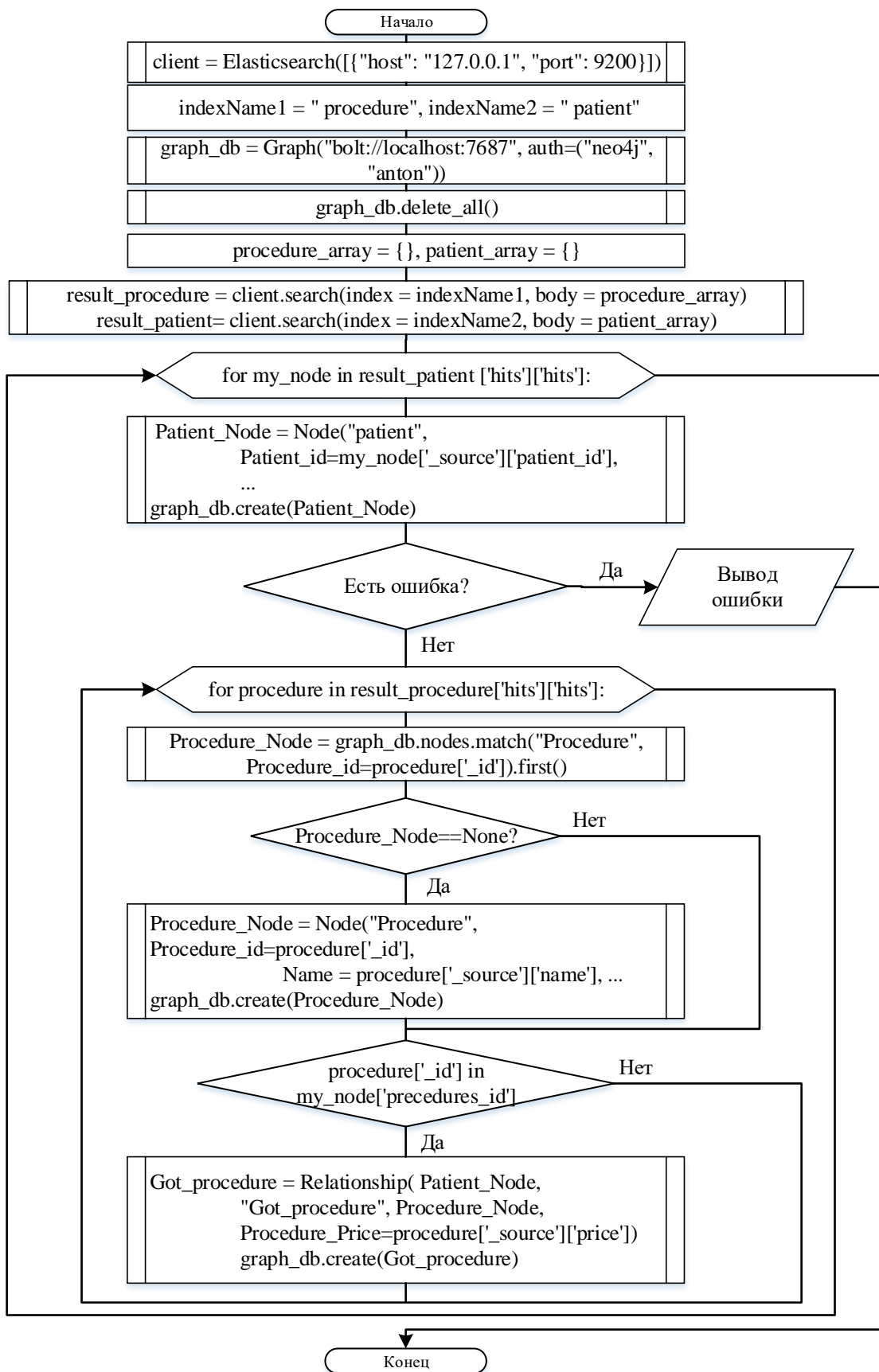
**Второй запрос:** определить стоимость процедуры по заданным ключевым словам.

```
1- {
2  "took" : 273,
3  "timed_out" : false,
4  "_shards" : {
5    "total" : 1,
6    "successful" : 1,
7    "skipped" : 0,
8    "failed" : 0
9  },
10 "hits" : {
11   "total" : {
12     "value" : 3,
13     "relation" : "eq"
14   },
15   "max_score" : 5.1222467,
16   "hits" : [
17     {
18       "index" : "procedure",
19       "type" : "Procedure",
20       "id" : "PR18",
21       "score" : 5.1222467,
22       "source" : {
23         "price" : "7000",
24         "name" : "Терапия травмами"
25       },
26       "highlight" : {
27         "description" : [
28           "Процедура, основанная на исполн>трав/ел> и растений для озд>укрепления"
29         ]
30       }
31     },
32     {
33       "index" : "procedure",
34       "type" : "Procedure",
35       "id" : "PR25",
36       "score" : 2.183346,
37       "source" : {
```

```
1 GET procedure/_search
2 {
3   "query": {
4     "simple_query_string": {
5       "query": "травы"
6     }
7   },
8   "_source": [
9     "name",
10    "price"
11  ],
12  "highlight": {
13    "fields": {
14      "description": {}
15    }
16  }
17 }
18 }
```

Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 6)

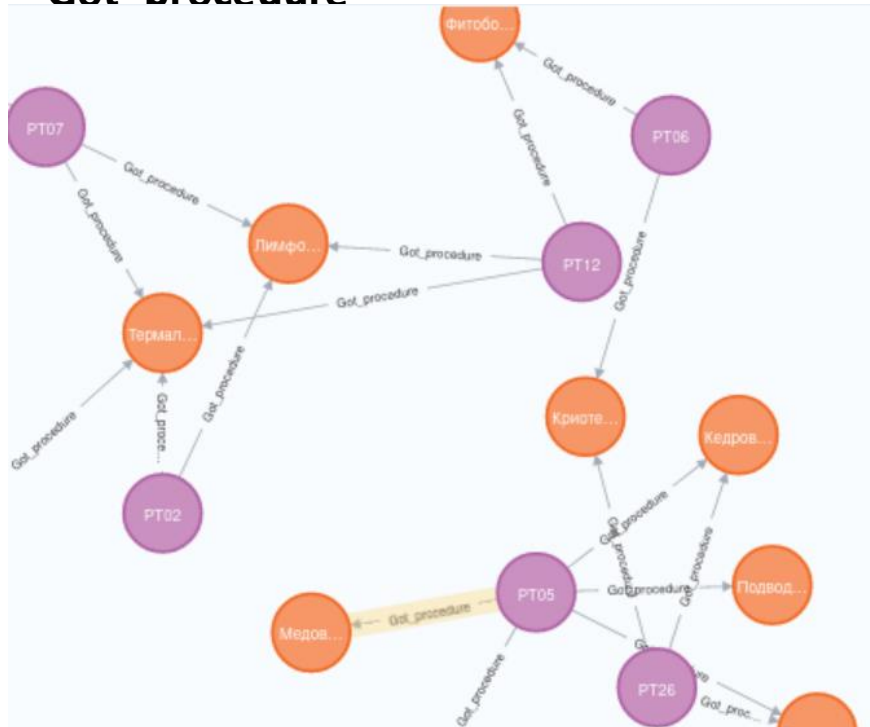
					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 6)											
					Elasticsearch. Запросы						Литер.		Масса	Масштаб		
Изм	Лист	№ документа	Подпись	Дата												
Разраб.																
Руковод.																
											Лист 4		Листов 10			
Н. Контр.																



					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 6)				
Изм	Лист	№ документа	Подпись	Дата	Нео4j. Алгоритм создания и заполнения графовой БД	Литер.	Масса	Масштаб	
Разраб.									
Руковод.									
Н. Контр.						Лист 5		Листов 10	

# Neo4j. Запрос

Связь «пациент получил процедуру» = Got procedure



Запрос: найти процедуру с максимальной суммарной стоимостью.

```
1 MATCH p=(pct:patient)-[r:Got_procedure]-(proc:Procedure)
2 WITH proc, sum(toInteger(r.Procedure_Price)) as procedure_sum, count(r) as number_of_client
3 ORDER BY procedure_sum desc
4 RETURN proc.Name, procedure_sum, proc.Price, number_of_client, proc
5 LIMIT 1
```

proc.Name	procedure_sum	proc.Price	number_of_client	proc
"Термальные ванны"	80000	"8000"	10	{ "identity": 1, "labels": [ "Procedure" ], "properties": { "Description": "Процедура, основанная на принятии термальных ванн с использованием теплой или горячей воды с добавлением различных минералов и микроэлементов, что способствует расслаблению"  } }

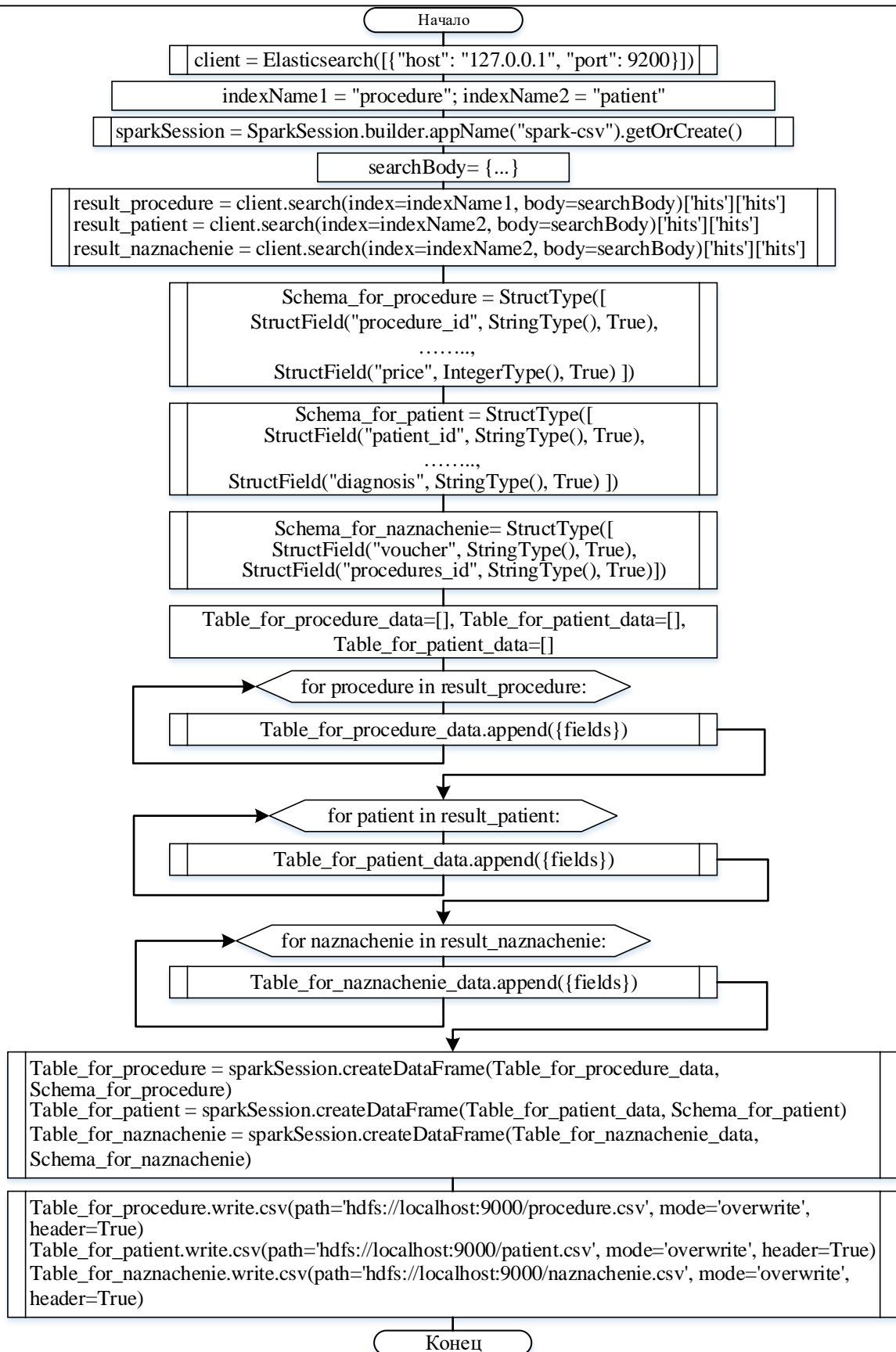
patient	
<id>	56
Arrival_date	2018-07-17
Diagnosis	[Депрессия,Панические атаки,Психосоматические расстройства]
Numb_of_days	11
Patient_id	PT09
Patient_personal_data	Смирнов Алексей Владимирович
Procedures_id	[PR01,PR03,PR30]
Voucher	4321567

## Node Properties

Procedure	
<id>	36
Description	Процедура, основанная на использовании лечебных трав и растений для оздоровления организма, укрепления иммунной системы и снятия стресса.
Name	Терапия травами
Price	7000
Procedure_id	PR18

Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)

					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)											
					Neo4j. Запрос						Литер.		Масса	Масштаб		
Изм	Лист	№ документа	Подпись	Дата												
Разраб.																
Руковод.																
											Лист 6		Листов 10			
Н. Контр.																



					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)				
Изм	Лист	№ документа	Подпись	Дата	Spark. Алгоритм создания CSV-файлов с таблицами	Литер.	Масса	Масштаб	
Разраб.									
Руковод.									
Н. Контр.									
						Лист 7		Листов 10	



# Spark. Запрос

Запрос: определить число выполненных процедур каждым специалистом

ID_OF_PROCEDURE	NUMBER_OF_NAZNACH
PR01	9
PR02	10
PR03	7
PR04	3
PR05	3
PR06	5
PR07	7
PR08	9
PR09	4
PR10	1
PR11	6
PR12	6
PR13	4
PR14	5
PR15	3
PR16	1
PR17	4
PR18	7
PR19	5
PR20	4
PR22	4
PR23	5
PR24	5
PR25	4
PR26	5
PR27	3
PR28	5
PR29	6
PR30	11

```
#!/usr/bin/env python
# -*- coding: utf-8 -*- from pyspark.sql import SparkSession

from pyspark.sql import *

sparkSession=SparkSession.builder.appName("Python Spark SQL basic
example").config("spark.sql.shuffle.partitions","100").getOrCreate()

Procedure_Table = sparkSession.read.load(path='hdfs://localhost:9000/procedure.csv', format='csv', sep=',',
inferSchema="true", header="true")
Master_Table = sparkSession.read.load(path='hdfs://localhost:9000/patient.csv', format='csv', sep=',',
inferSchema="true", header="true")

Procedure_Table.registerTempTable("procedure")
Patient_Table.registerTempTable("patient")

df = sparkSession.sql("SELECT ID_OF_PROCEDURE, COUNT(ID_OF_PROCEDURE) as NUMBER_OF_NAZNACH
FROM (SELECT procedure.procedure_id as ID_OF_PROCEDURE, naznachenie.procedures_id LIKE concat(concat('%',
procedure.procedure_id),'%') as CHECK, FROM procedure, naznachenie)
WHERE check=true
GROUP BY ID_OF_PROCEDURE
ORDER BY ID_OF_PROCEDURE").show(30)

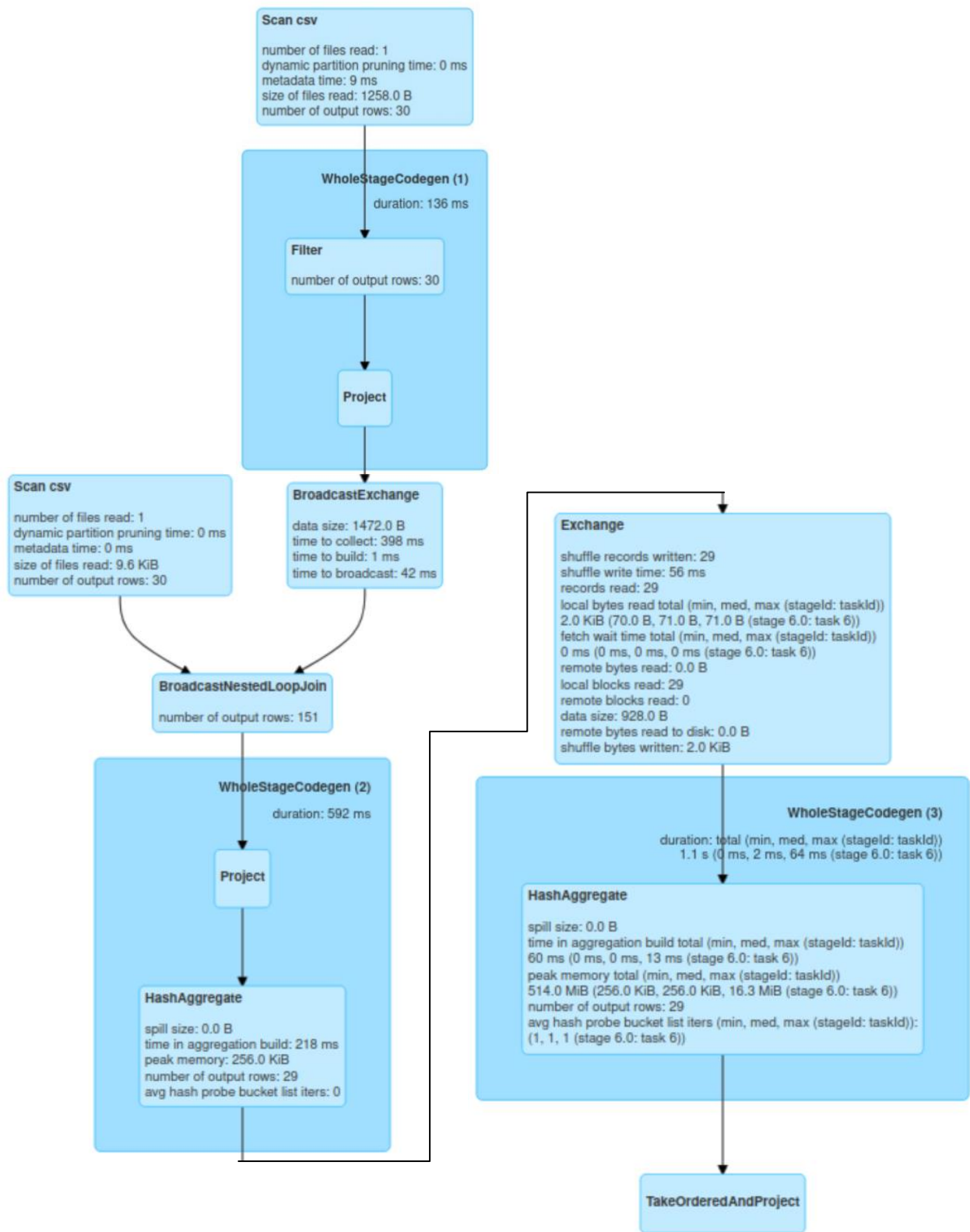
input('Ctrl C')
```

					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)				
					Spark. Запрос		Литер.	Масса	Масштаб
Изм	Лист	№ документа	Подпись	Дата					
Разраб.									
Руковод.									
							Лист 8		Листов 10
Н. Контр.									



					Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)						
					Spark. Мониторинг выполнения запроса	Литер.			Масса	Масштаб	
Изм	Лист	№ документа	Подпись	Дата							
Разраб.											
Руковод.											
						Лист 9			Листов 10		
Н. Контр.											

# Spark. DAG SQL-запроса



Разработка макета аналитической системы на основе баз данных NoSQL (вариант № 24)

Изм	Лист	№ документа	Подпись	Дата	Spark. DAG SQL-запроса	Литер.	Масса	Масштаб
Разраб.								
Руковод.								
Н. Контр.						Лист 10	Листов 10	