



**EuroHPC**  
Joint Undertaking

# **„Възможности за интегриране на Nadoor в RapidMiner за анализ на големи данни“**

## **Съдържание**

<b>Въведение .....</b>	<b>2</b>
<b>Проучване и анализ на технологичния проблем .....</b>	<b>2</b>
<b>Приложното решение .....</b>	<b>9</b>
<b>Обработка на големи данни .....</b>	<b>16</b>
<b>Използвана литература .....</b>	<b>21</b>

## Въведение

Големите данни предоставят възможности за трансформиране на организациите като засягат различни аспекти на организационните дейности. Анализът на големите данни, генерирани в рамките на ежедневните бизнес дейности в днешно време, е чудесна възможност за организациите да подобрят своите процеси на вземане на решения и стратегическо управление. Организациите, които използват големи данни, могат да отключат множество възможности за получаване на ценна информация, подобряване на процесите на вземане на решения и цялостната производителност. Големите данни изискват инфраструктура, която да осигури възможности за обработка и анализ на данните.

Hadoop е рамка с отворен код, предназначена да позволи разпределената обработка на големи набори от данни в клъстер от компютри, осигурявайки толерантност към грешки и мащабируемост. Със своята разпределена файлова система (HDFS) и възможности за паралелна обработка (MapReduce), Hadoop се използва широко за съхраняване и анализиране на огромни количества структурирани, полуструктурирани и неструктурирани данни.

RapidMiner, от друга страна, е интегрирана среда за машинно обучение, извличане на данни, прогнозни анализи и разширен анализ. RapidMiner улеснява откриването и използването на модели и връзки в данните, позволявайки на организациите да вземат информирани решения и да извличат полезно знание. Разполага с удобен за потребителя интерфейс и голям набор от инструменти. Проектиран е да обработва различни типове данни и размери, включително големи данни. RapidMiner предоставя няколко начина за работа с големи данни, което го прави подходящ за широкомащабни задачи за анализ на данни.

Интегрирането на Hadoop в RapidMiner позволява разпределена обработка и съхранение на големи данни от Hadoop, и възможности, предлагани от RapidMiner, за ефективен анализ на данните и прилагане на машинно обучение.

## Проучване и анализ на технологичния проблем

Инфраструктурата, необходима за анализа на големи данни, обикновено включва комбинация от хардуер, софтуер и мрежови компоненти. Тъй като големите данни включват обработка и съхранение на огромни обеми от данни, инфраструктурата трябва да бъде мащабируема, надеждна и способна да обработва сложни задачи за обработка на данни ефективно.

Основни компоненти на инфраструктурата, необходими за анализ на големи данни са:

- **Системи за съхранение на данни**
  - Разпределени файлови системи: Разпределена файлова система на Hadoop (HDFS) и Apache Hadoop съвместима файлова система (HCFS) позволяват данните да се разпространяват и съхраняват в множество възли в клъстер.
  - NoSQL бази данни: Apache Cassandra, MongoDB или Apache HBase, които са подходящи за обработка на неструктурирани и полуструктурирани данни.
- **Изчислителни ресурси**
  - Клъстерни изчисления: Клъстер от взаимосвързани сървъри или възли за разпределяне на задачи за обработка на данни и паралелно изпълнение на изчисления.
  - Високопроизводителни изчисления (HPC): Мощни сървъри или възли с множество ядра и достатъчно памет, за да се справят с ресурсоемки задачи за анализ на данни.
- **Поглъщане и интегриране на данни**
  - Интегриране на данни: Инструменти за поглъщане, изчистване и интегриране на данни от различни източници в системите за съхранение на данни.
  - Потоци за данни: За оркестриране и автоматизиране на работните потоци от данни от поглъщане до анализ.
- **Обработка на данни и анализи**
  - Обработка на големи данни: Apache Hadoop MapReduce, Apache Spark или Apache Flink, за разпределена обработка на данни в клъстера.
  - Аналитични инструменти: Apache Hive, Apache Pig или Apache Impala, за търсене и анализиране на данни с помощта на езици, подобни на SQL.

- Библиотеки за машинно обучение: За извършване на усъвършенствани анализи, прогнозно моделиране и задачи за машинно обучение на големи данни.
- **Визуализация на данни**
  - Инструменти за визуализация на данни с цел създаване на интерактивни и информативни визуални представяния на анализирани данни.
  - Платформи за управленски табла, които да генерират изчерпателни отчети и табла за вземане на решения.
- **Мониторинг и управление**
  - Управление на клъстери за мониторинг на ефективността на инфраструктурата и на клъстерните възли.
  - Системи за управление на ресурсите: Apache YARN или Kubernetes за ефективно разпределяне и управление на изчислителните ресурси в клъстера.
- **Сигурност и управление**
  - Мерки за осигуряване на сигурност на данните като се гарантира поверителността на данните и да се защити чувствителната информация.
  - Механизми за контрол на достъпа с цел управление на потребителските разрешения и правата за достъп до данни и ресурси.
  - Политики за управление на данни чрез установяване на правила и насоки за управление и използване на данни.
- **Мрежова инфраструктура**
  - Високоскоростни мрежи с цел да се осигури бърз трансфер на данни между възлите и ефективна комуникация в клъстера.
- **Облачни услуги:**
  - Инфраструктура, базирана на облак: Доставчиците на публични облаци предлагат мащабируема и гъвкава инфраструктура за анализ на големи данни, без да е необходимо локално управление на хардуера.

Специфичните изисквания за инфраструктурата могат да варират в зависимост от обема на данните, сложността на анализа и конкретните случаи на използване на големите данни. Организацията могат да изберат да изградят и поддържат своята инфраструктура локално, да използват облачни услуги или да възприемат хибриден подход въз основа на техните нужди и ресурси.

HDFS (Hadoop Distributed File System) е разпределена файлова система, предназначена да съхранява и управлява големи обеми от данни в множество машини в Hadoop клъстер. HDFS е подходящ за съхранение и управление на големи данни, структурирани, полуструктурирани и неструктурирани данни, което го прави основен компонент на много приложения и работни потоци за големи данни.

HDFS разпределя данни в множество възли в клъстера Hadoop. Той разбива големите файлове на по-малки блокове, обикновено с размер по подразбиране от 128 MB или 256 MB, и съхранява множество копия на всеки блок на различни възли, за да осигури толерантност към грешки. HDFS е проектиран да бъде толерантен към грешки. Ако възел или диск се повреди, системата може автоматично да възпроизведе данните, съхранени на този възел, в други здрави възли, осигурявайки наличност на данни дори при наличие на хардуерни повреди.

HDFS е оптимизиран за задачи за пакетна обработка на данни, но осигурява достъп до данни с висока пропускателна способност за четене и писане на данни по поточен начин.

HDFS следва WORM (Write-Once-Read-Many) модел, което означава, че данните могат да бъдат записани във файловата система веднъж и след това стават само за четене. Това е подходящо за обработка на големи данни, при които данните рядко се актуализират, след като се съхраняват. HDFS репликира блокове от данни в множество възли, за да осигури излишък на данни. Коефициентът на репликация по подразбиране обикновено е три, което означава, че всеки блок данни се съхранява на три различни възела. Когато обработва данни, рамката на Hadoop се опитва да планира задачи на възли, където се намират данните, като минимизира трансфера на данни в мрежата.

Apache Hive е език за съхранение на данни и SQL-подобен език за заявки, който е изграден върху екосистемата на Hadoop. Той осигурява слой над разпределените възможности за съхранение и обработка на Hadoop, което улеснява потребителите да взаимодействат и анализират големи набори от данни, използвайки познатия SQL синтаксис.

Основните характеристики на Apache Hive включват:

- Hive използва SQL-подобен език за заявки, наречен HiveQL (Hive Query Language) или HQL. Потребителите могат да пишат заявки с помощта на SQL синтаксис, което улеснява извършването на анализ на големи данни, съхранявани в Hadoop.
- Hive е схема-на-четене, което означава, че позволява на потребителите да наложат структура на данните, когато ги четат от основното хранилище Hadoop (напр. HDFS). Това позволява на потребителите да правят заявки и да анализират данни, без да изискват предварително дефинирана схема по време на поглъщане на данни.
- Hive е оптимизиран за пакетна обработка и е подходящ за задачи за обработка на данни. Въпреки че предоставя функционалност, подобна на SQL, тя може да не е подходяща за заявки в реално време или интерактивни заявки.
- Hive се интегрира с други компоненти на екосистемата на Hadoop, като HDFS за разпределено съхранение и YARN за управление на ресурсите. Той може да работи и с данни, съхранявани в HBase, NoSQL база данни, изградена върху Hadoop.
- Поддържа разделяне и групиране на данни, което позволява на потребителите да организират данните по-ефективно за по-бързо търсене и обработка.
- Hive поддържа дефинирани от потребителя функции (UDFs) и дефинирани от потребителя агрегати (UDAs), което позволява разширяване на функционалността чрез писане на персонализирани функции на Java, Python или други езици за програмиране.
- Hive използва metastore за съхраняване на информация за метаданни за таблици, дялове и данни, съхранявани в Hadoop.

Докато Hive опростява анализа на големи данни, използвайки синтаксис, подобен на SQL, той може да не е толкова ефективен за заявки в реално време или интерактивни заявки. За тези случаи на употреба, други инструменти като Apache

Impala или Apache Spark SQL може да са по-подходящи. Hive е най-подходящ за пакетна обработка и аналитични задачи на големи набори от данни в Hadoop.

Hive може да се интегрира с външни системи, като Apache HBase или решения за съхранение в облак, за достъп и анализ на данни, съхранявани в тези системи.

RapidMiner е платформа за анализ на данни, която включва усъвършенствани възможности за анализ и машинно обучение. Той е проектиран да обработва различни типове и размери данни, включително големи данни.

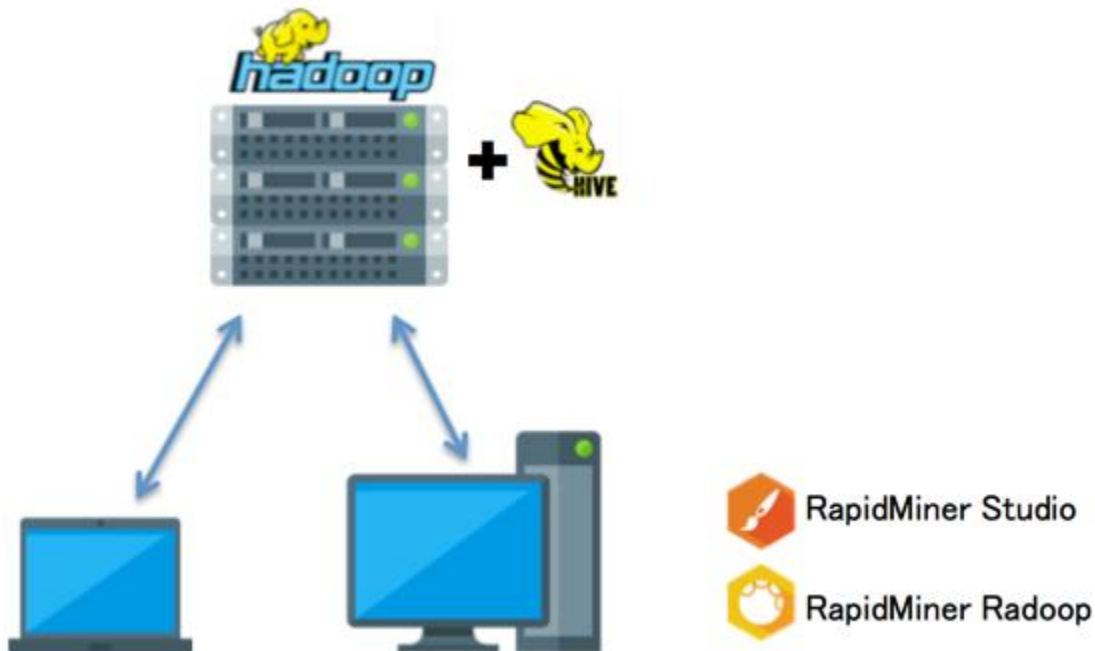
RapidMiner предоставя няколко начина за работа с големи данни, което го прави подходящ за задачи за анализ на големи данни.

- **Конектори за големи данни:** RapidMiner предлага конектори към различни системи за съхранение на големи данни, като Hadoop разпределена файлова система (HDFS), Apache Hive, Apache Spark и Apache HBase. Тези конектори позволяват директен достъп и четене на данни, съхранявани в тези разпределени системи за съхранение за анализ.
- **Разпределени изчисления с Spark:** RapidMiner позволява да се използва Apache Spark за разпределена обработка на данни. С помощта на Spark може да се извършват задачи за трансформации на данни, анализи и машинно обучение на големи набори от данни по паралелен и разпределен начин.
- **Обработка в базата данни:** RapidMiner поддържа обработка в базата данни, където обработката и анализът на данни се изтласкват надолу до ниво база данни.
- **Извадки от данни и филтриране на данни:** За да се справи с големи набори от данни, RapidMiner предоставя различни опции за формиране на извадка и филтриране на данни.
- **Паралелна обработка и оптимизация:** RapidMiner оптимизира обработката на данни чрез използване на техники за паралелизация, с цел максимално използване на хардуерните ресурси и намаляване времето, необходимо за анализ.

- **Мащабируемост и производителност:** RapidMiner обработва ефективно големи данни и може да се мащабира, за да се приспособи към нарастващите набори от данни и повишените изисквания за обработка.
- **Предварителна обработка на данни:** RapidMiner предлага широка гама от оператори за предварителна обработка на данни за почистване, трансформиране и подготвяне на големи данни за анализ.
- **Машинно обучение и прогнозно моделиране:** Възможностите за машинно обучение на RapidMiner могат да бъдат приложени към големи данни, което позволява изграждане и внедряване на аналитични модели на големи данни.

RapidMiner Radoop е продукт на RapidMiner, който осигурява графичен интерфейс за анализ на големи данни в Hadoop клъстер с работещ Hive сървър. Radoop изисква Hadoop клъстера да бъде достъпен от клиента, изпълняващ RapidMiner Studio.

Диаграмата по-долу показва основната архитектура на решението RapidMiner Radoop на RapidMiner Studio:





Фигура 1 Архитектура на решението RapidMiner Radoop на RapidMiner Studio<sup>1</sup>

Функциите на RapidMiner Radoop се инсталират автоматично в базата данни Hive, конфигурирана във връзката в Radoop.

RapidMiner Radoop автоматично качва два файла (radoop\_hive-vX.jar и rapidminer\_libs-<version>.jar) в HDFS и ги използва за дефиниране на персонализирани Hive функции (UDFs). За защитени Hadoop клъстери това може да е забранено, така че администраторът на Hadoop трябва да инсталира тези UDF. Има два начина за това:

- при Cloudera Distribution, се използва Parcel, предоставен от RapidMiner
- Ръчно инсталиране на jar файлове

### Приложното решение

Интеграцията на Hadoop и RapidMiner предлага множество подходи в зависимост от специфични нужди [1].

- **Директна Hadoop интеграция:** RapidMiner осигурява директна интеграция с Hadoop, което позволява директно взаимодействие и обмен на данни с разпределената файлова система и ресурсния мениджър на Hadoop. Този интеграционен подход предлага безпроблемна свързаност и оптимизирано използване на разпределените възможности за обработка на Hadoop в рамките на RapidMiner.
- **Интеграция, базирана на конектори:** RapidMiner предлага конектори към различни съвместими с Hadoop системи за съхранение, като HDFS, Hive, Impala, или Apache HBase, което позволява взаимодействие с Hadoop данни чрез тези конектори, извличайки, трансформирайки и зареждайки данни в RapidMiner без директна интеграция.
- **Интеграция чрез командния ред:** За по-напреднали потребители или специфични случаи на употреба, RapidMiner предоставя интерфейс, който позволява директно изпълнение на процесите на RapidMiner в клъстера Hadoop. Този подход позволява използването на разпределените възможности

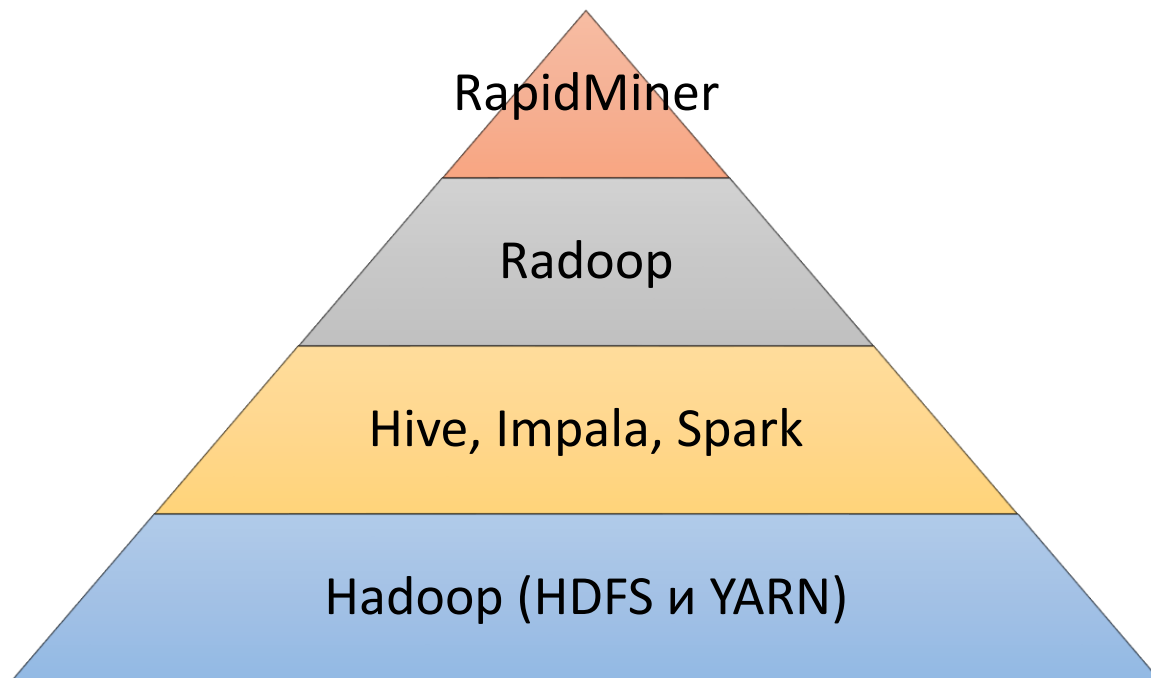
---

<sup>1</sup> <https://docs.rapidminer.com/7.6/radoop/overview/>

за обработка на Hadoop, като същевременно се използват аналитичните функционалности на RapidMiner.

- **Интеграция с други рамки за големи данни:** Освен Hadoop, RapidMiner се интегрира и с други популярни рамки за големи данни като Apache Spark или Apache Flink. Тези рамки осигуряват възможности за обработка в реално време и мащабируема обработка на данни, отваряйки допълнителни възможности за анализи с RapidMiner.

RapidMiner Radoop се инсталира като разширение ц RapidMiner Studio. Конфигурирането на връзка между RapidMiner Radoop в RapidMiner Studio и един или повече Hadoop клъстери от Управление на връзките на Radoop и Настройки на връзката. Достъп до тези връзки има от менюто Connections, изгледа Hadoop Data или изгледа Design. След като секонфигурира и запази връзката, може да се тества преди създаване на процесите на обработка и анализ. Тестът валидира връзката към клъстера и проверява дали настройките на връзката отговарят на изискванията на RapidMiner Radoop.




Изискванията на RapidMiner Radoop към Hadoop клъстера и хранилището са:

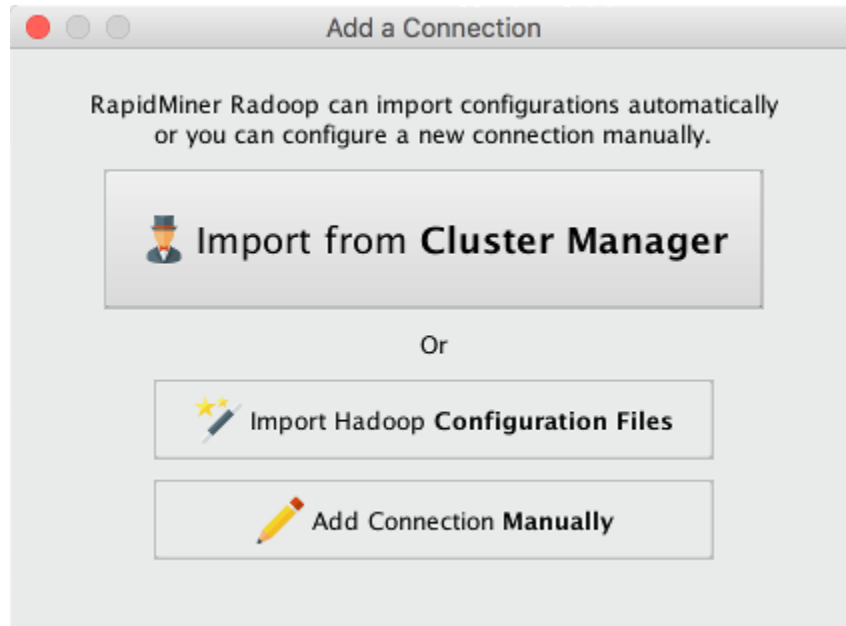
- поддържана дистрибуция на Hadoop, която се състои от **HDFS и YARN**.



Поддържаните дистрибуции са:

- Amazon Elastic MapReduce (EMR) 4.4+
  - Apache Hadoop 2.2+
  - Apache HDInsight 3.5
  - Cloudera Hadoop CDH5.x
  - Hortonworks HDP 2.x
  - IBM Open Platform 4.1+
  - MapR 5.x
  - Open Data Platform 0.9+
- разпределена система за съхранение на данни (**Hive или Impala**). RapidMiner Radoop поддържа следните инфраструктури за съхранение на данни:
    - Apache HiveServer2 0.13+
    - Cloudera Impala 1.2.3 and later (see Impala limitations on the Installing Radoop on Studio page)
  - **Java 8** инсталиран на клъстерните възли (необходими за прилагане на повечето модели на RapidMiner в Hadoop и използване на оператори за процеси). RapidMiner Radoop изисква Java 8, инсталиран на клъстера Hadoop, за да работи. Възлите трябва да имат поне 8 GB RAM.
  - **Apache Spark** (по желание). RapidMiner Radoop поддържа следните версии на Spark:
    - Apache Spark 1.2.x, 1.3.x and 1.4.x. Поддържа оператори за дърво на решенията, линейна регресия и логистична регресия.
    - Apache Spark 1.5.x, 1.6.x, 2.0.x (освен 2.0.1), 2.1.x, 2.2.x. Поддържа всички оператори на Spark, включително Spark Script (Python и / или R се изисква на клъстерните възли), Single Process Pushdown и SparkRM.
    - Версия на Spark 2.0.1 не се поддържа.

Три метода за създаване на Radoop връзка:

- Ако има достъп до **софтуер** за управление на клъстери (Apache Ambari или Cloudera Manager), препоръчва се използването на  **Import from Cluster Manager**. Този метод е най-лесен.



- Ако не се използва или липсва достъп до клъстерен мениджър, може да се използват *конфигурационните файлове на клиента* в  **Import Hadoop Configuration Files**.
- В противен случай има възможност за ръчна конфигурация  **Add Connection Manually**.

Настройките на връзката има няколко раздела.

#### Глобални настройки:

- Версия на Hadoop - дистрибуцията, която дефинира типа версия на Hadoop за тази връзка.
- Допълнителна директория за библиотеки - всички допълнителни библиотеки (JAR файлове) на клиента, необходими за свързване към клъстера
- Настройки за връзка с Hadoop клъстер, защитен от Kerberos:
  - Client Principal - Потребителят, който има достъп до Hadoop. Форматът е `primary[/<instance>]@<REALM>`, where primary is usually the user name,

instance is optional, and REALM is the Kerberos realm. Пример: [user/client.rapidminer.com@RAPIDMINER.COM](mailto:user/client.rapidminer.com@RAPIDMINER.COM).

- Използване на парола вместо keytab file – удостоверяване с парола вместо с KeyTab File.
- KeyTab File - път към keytab file на клиентската машина. Въведете или намерете местоположението на файла.
- Парола - паролата на Kerberos, която може да се използва за свързване със защитения клъстер. RapidMiner Radoop използва файла .key шифър, за да криптира паролата в radoop\_connections.xml.
- KDC Address - адрес на Kerberos Key Distribution Center. Пример: kdc.rapidminer.com.
- REALM - обикновено това е името на домейна с главни букви. Пример: RAPIDMINER.COM.
- Kerberos Config File - за да се избегнат конфигурационни разлики между машината, работеща с RapidMiner и клъстера Hadoop, добра практика е да предоставите конфигурационния файл Kerberos (обикновено krb5.conf или krb5.ini). Файлът се получава от администратора. Въведете или намерете местоположението на файла.
- Настройки на Hadoop MapR
  - Enable MapR security - свързване с Hadoop клъстер, защитен от MapR Security.
  - MapR cluster - MapR клъстер за свързване. Всички MapR връзки трябва да бъдат конфигурирани в MapR клиента, към който сочи MapR Home.
- Hadoop Username - името на потребителя на Hadoop. В повечето случаи потребителят трябва да има подходящи разрешения за клъстера. За нова връзка по подразбиране е потребителят на операционната система.

## Hadoop:

### Настройки на NameNode:

- Адрес (обикновено име на хост) на възела, изпълняващ услугата NameNode. (Изисква работеща система за разрешаване на имена на мрежи.)
- Порт на услугата NameNode.
- Адрес (обикновено име на хост) на възела, изпълняващ услугата "Диспечер на ресурси".
- Порт на услугата "Диспечер на ресурси".
- Адрес (обикновено име на хост) на възела, изпълняващ услугата Job History Server.
- Порт на услугата Job History Server.
- Настройки при активирана защита на Kerberos и деактивирано извличане на Hive.
  - Retrieve Service Principals от Hive - Ако е отметнато, RapidMiner Radoop автоматично извлича всички други директори на услуги от Hive за по-лесно конфигуриране. Забраняване на тази настройка само ако има проблем с достъпа до други услуги.
  - NameNode Principal- може да се използва ключовата дума \_HOST като пример. Пример: [nn/\\_HOST@RAPIDMINER.COM](mailto:nn/_HOST@RAPIDMINER.COM)
  - Resource Manager Principal - може да се използва ключовата дума \_HOST като пример. Пример: [rm/\\_HOST@RAPIDMINER.COM](mailto:rm/_HOST@RAPIDMINER.COM)
  - JobHistory Server Principal - може да се използва ключовата дума \_HOST като пример. Пример: [jhs/\\_HOST@RAPIDMINER.COM](mailto:jhs/_HOST@RAPIDMINER.COM)
- Разширени параметри на Hadoop - свойства на ключовата стойност за персонализиране на връзката Hadoop и задачите на Radoop Yarn / MapReduce . Някои връзки изискват определени разширени параметри.

## Spark

- Spark версия инсталирана на клъстера.
- Assembly Jar местоположение/Spark Archive (или libs) път - местоположението на HDFS или локалния път (на всички клъстерни възли) на файла Spark Assembly Jar / Spark Jar файлове.
- Правила за разпределение на ресурси на Spark – политиката за разпределение на ресурси за Spark задания.
- Разпределение на ресурси % - Процент от ресурсите на клъстера, разпределени за задание на Spark. Това поле е разрешено само когато "Статична, евристична конфигурация" е политиката за разпределение на ресурсите на Spark.
- Разширени параметри - свойства на ключовата стойност, които персонализират заданията на Spark на RapidMiner Radoop.

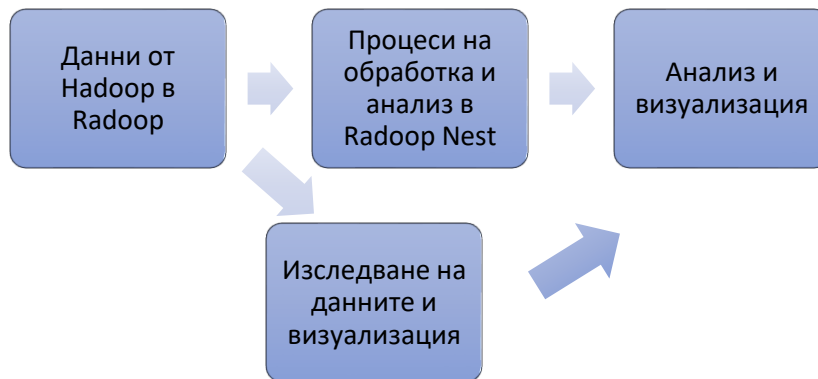
## Hive

- Hive Версия - подходящата система за съхранение на данни — HiveServer2 (Hive 0.13 или по-нова) или Impala. Като алтернатива може да се избере Custom HiveServer2. В случай че е избран Custom HiveServer2 се избира директория, която съдържа библиотеките (JAR файлове), необходими за свързване към клъстера.
- Hive Server Address/Impala Server Address
- Hive Port/Impala Port
- Database name
- Username - Потребителско име за свързване към указаната база данни. По подразбиране е "hive" за всички връзки с версии на HiveServer2. Този потребител трябва да има достъп до директорията HDFS, която Radoop използва за временно съхраняване на файлове. Ако тази директория се намира в зона за шифроване, потребителят също трябва да има разрешения за достъп до ключа на зоната за шифроване.

- Парола за свързване към указаната база данни. RapidMiner Radoop използва шифъра.key файла, за да криптира паролата в radoop\_connections.xml.

## Обработка на големи данни

Основните стъпки за обработка на големи данни в Radoop са представени на фигурата по-долу.



След зареждането в Radoop от Hadoop Distributed File System (HDFS), HIVE таблици и Apache Spark DataFrames, данните могат да се изследват в Radoop чрез използването на различни оператори за изследване и анализ.

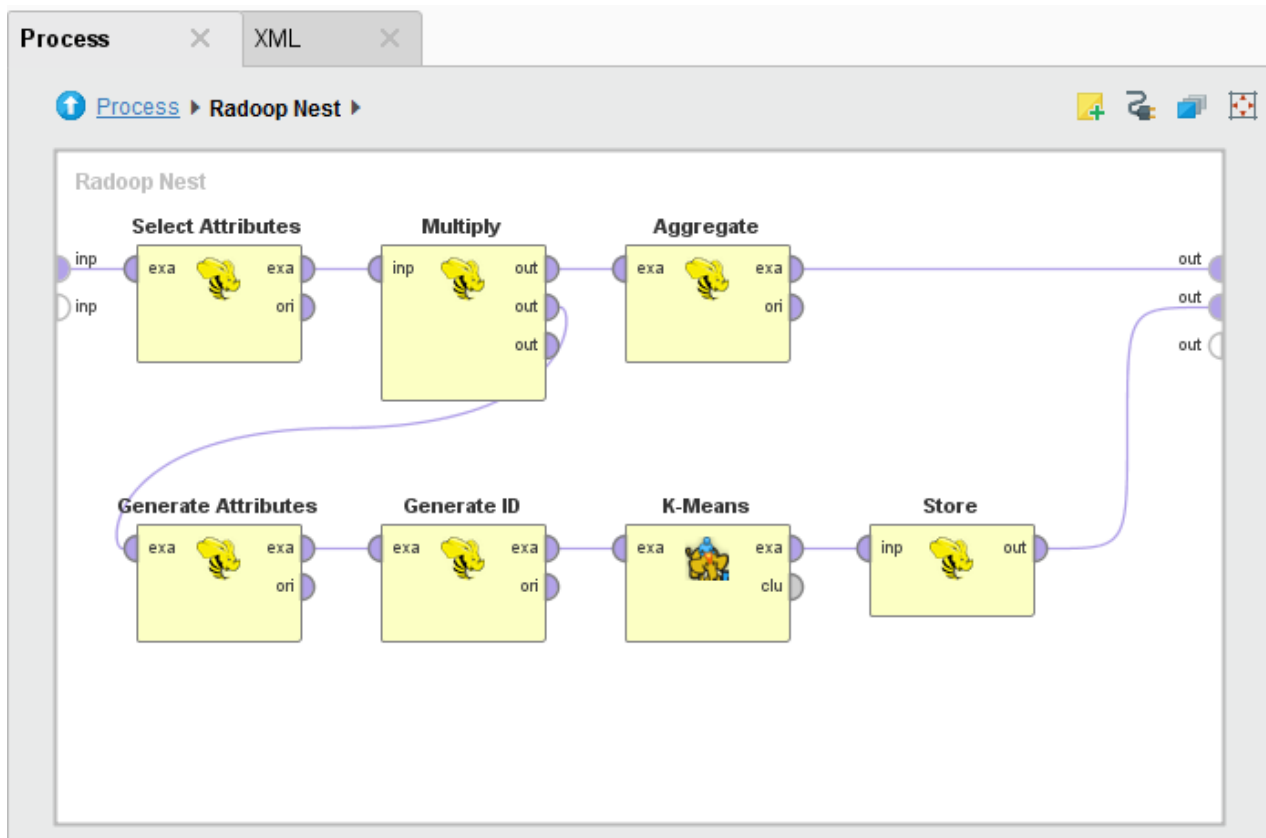
- **Трансформация на данни:** След като данните се заредят, може да се извършат различни операции за трансформация на данни, като се използват разпределените изчислителни възможности на Radoop. Radoop осигурява широк спектър от функции за манипулиране на данни, подобни на наличните в R.
- **Изследване и анализ на данните:** С Radoop може да се извършва разпределено извличане на данни, машинно обучение и статистически анализи на големи масиви от данни.
- **Изграждане на модели:** Radoop поддържа изграждане на модели за машинно обучение, използвайки разпределени алгоритми, предоставени от Apache



Spark. Radoop се интегрира с разпределените библиотеки за машинно обучение на Apache Spark за изграждане и обучение на модели за машинно обучение на големи данни.

- **Оценка на модела:** След изграждането на моделите, може да се направи оценка на тяхното представяне чрез съответните показатели за оценка.
- **Внедряване и интеграция:** След като анализът приключи, моделът може да се внедри или да се интегрира в съществуващи потоци за обработка на данни.
- **Наблюдение и оптимизиране:** Непрекъснатото наблюдение на обработката на данни и производителността на модела.

На фигурата по-долу е представен примерен процес в Radoop Nest за обработка и анализ на големи данни.



Фигура 2 Примерен процес в Radoop Nest за обработка и анализ на големи данни.

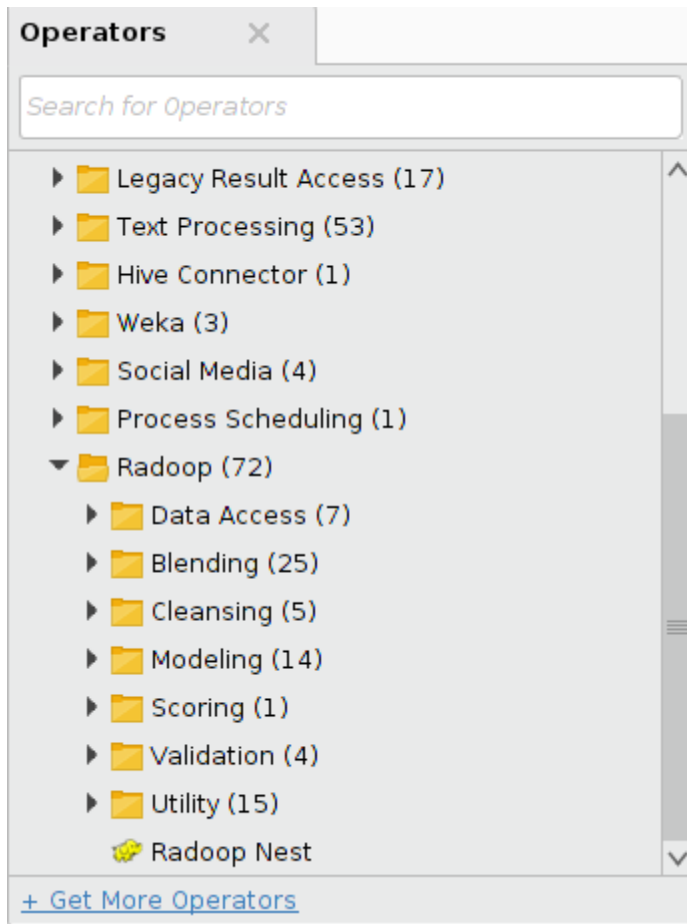
Radoop Nest е най-важният елемент в RapidMiner Radoop. Всеки процес трябва да съдържа поне един Radoop Nest (мета) оператор; той определя връзката с Hadoop клъстера. Всеки подпроцес на операторите на Radoop, поставени в Radoop Nest, описва процеса, който се изпълнява на този клъстер; цялото съдържание извън Radoop Nest (обработено от операторите на RapidMiner Studio) се обработва в паметта. Повечето оператори на Radoop имат аналог в RapidMiner Studio. Основната разлика между двете е, че операторите на Radoop винаги съхраняват и обработват данни в клъстера.

Radoop позволява да се комбинират базирани на памет и клъстерни оператори в един и същ процес. На входния си порт Radoop Nest импортира данни от оперативната памет на клиента в клъстера. Операторите вътре в Nest консумират и произвеждат обекти HadoopExampleSet (базираният на клъстери вариант на стандартния обект ExampleSet). HadoopExampleSet съхранява данните в Hive, във временна или постоянна таблица.

Radoop Nest може да има произволен брой изходни портове, за да достави базираните на памет ExampleSet обекти директно към изходен порт на процеса или към входния порт на следващия оператор на RapidMiner извън Nest. Radoop извлича данните или извадка от данни от изхода на HadoopExampleSet в оперативната памет на клиента и след това се разпространява по-нататък като базиран на паметта ExampleSet на потока на процеса. Тъй като извадката от данни трябва да се побере в оперативната памет, може да се анализират обобщени данни (след като агрегирането е извършено в клъстера).

В допълнение към Radoop Nest, който е контейнерът за подпроцеса на Hadoop клъстера, има много оператори на RapidMiner Radoop, които могат да се използват. Тези оператори са категоризирани в следните групи:

- Data Access
- Blending
- Cleansing
- Modeling
- Scoring
- Validation
- Utility



Radoop може да се използва в различни индустрии за обработка и анализ на големи данни:

- Предсказване отлив на клиенти: Radoop може да анализира данните за поведението на клиентите, за да предскаже вероятността от напускане. Чрез идентифициране на клиентите, изложени на риск от напускане, организациите могат да предприемат проактивни мерки, за да ги задържат.
- Системи за препоръки: С Radoop може да се разработват системи за препоръки, които предоставят персонализирани препоръки за продукти или съдържание въз основа на поведението и предпочитанията на потребителите.
- Персонализиране на офертите в електронната търговия: Онлайн компания за търговия може да използва Radoop, за да анализира данните за поведението на клиентите, историята на покупките и взаимодействията с уебсайтове, за да

предостави персонализирани препоръки за продукти, насочени маркетингови кампании и да подобри опита на клиентите.

- Оптимизация на телекомуникационната мрежа: Телекомуникационните компании могат да използват Radoor, за да анализират огромни обеми мрежови данни, да идентифицират мрежовите проблеми, да оптимизират разпределението на ресурсите и да прогнозират потенциални неуспехи, за да подобрят производителността на мрежата и да намалят времето на престой.
- Откриване на измами в здравеопазването: Доставчиците на здравни услуги и застрахователните компании могат да използват Radoor, за да анализират данните за здравни претенции, досиетата на пациентите и информацията за фактуриране, за да открият измамни дейности и да предотвратят неправомерни плащания.
- Производствен контрол на качеството: Производствените фирми могат да използват Radoor, за да анализират сензорните данни от производствените процеси, за да прогнозират и предотвратят проблеми с качеството, да оптимизират ефективността на производството и да намалят дефектите на продукта.
- Анализ на настроеността в социалните медии: Компаниите могат да използват Radoor, за да анализират данните от социалните медии, обществените настроения и обратната връзка с клиентите, за да получат представа за възприемането на марката, да идентифицират нововъзникващите тенденции и да отговорят проактивно на изискванията на клиентите.
- Оценка на финансовия риск: Финансовите институции могат да прилагат Radoor, за да анализират финансови данни, кредитна история и икономически показатели, за да оценят кредитния риск, да открият потенциални измами и да вземат решения за кредитиране, базирани на данни.
- Оптимизация на потреблението на енергия: Компаниите за комунални услуги могат да използват Radoor за обработка на данни от интелигентни измервателни уреди и исторически модели на потребление на енергия, за да

оптимизират разпределението на енергията, да идентифицират възможностите за пестене на енергия и да управляват ефективно пиковото потребление.

- Планиране на транспорта и логистиката: Транспортните компании могат да анализират големи обеми от данни с цел анализ на трафика, движението на превозните средства и маршрутите за доставка и да оптимизират логистичните операции, да намалят времето за доставка и да подобрят планирането на маршрута.

### Използвана литература

1. How to Integrate Hadoop Into RapidMiner, <https://zebrabi.com/advanced-guide/how-to-integrate-hadoop-into-rapidminer/> , July 21, 2023
2. Rapidminer Documentation, <https://docs.rapidminer.com/7.6/radoop/overview/>
3. Identify and Prevent Fraud, <https://rapidminer.com/customer-stories/identify-prevent-fraud/>