

Párhuzamos R programok végrehajtása sparklyr segítségével az ELKH Cloud infrastruktúráján

Understanding the role of Apache Spark in the big data ecosystem

Big Data

The theoretical approach

- a) Volume
- b) Variety
- c) Velocity

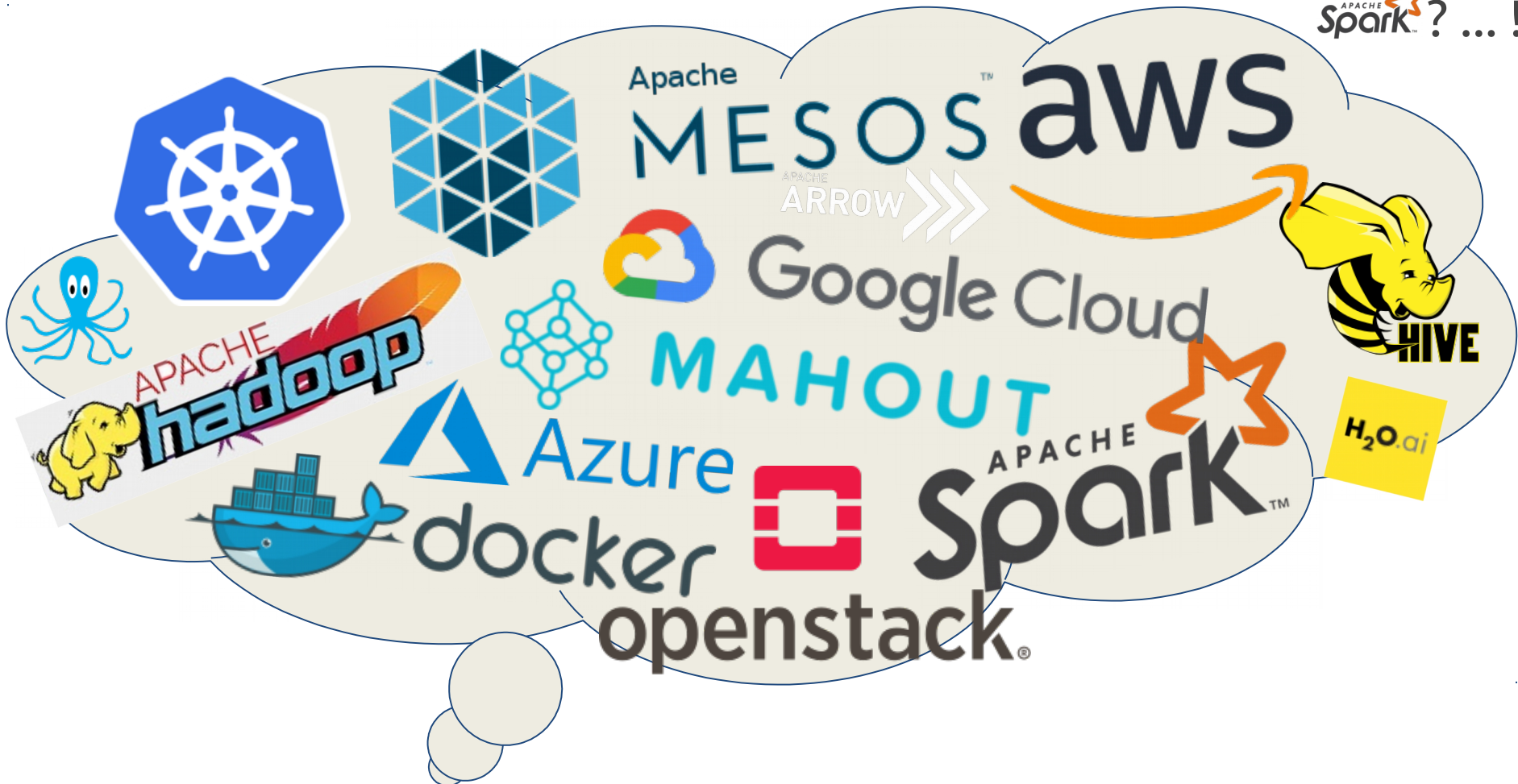
Big Data

The theoretical approach

- a) Volume
- b) Variety
- c) Velocity

The practical approach

- d) Doesn't fit
- e) Too slow



It can feel a bit overwhelming. (And this is just the tip of the iceberg.)

Implementing a machine learning process: the **how**

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations is achieved

Data distribution and **I/O operations**

Implementing a machine learning process: the **how**

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations is achieved

Data distribution and I/O operations



Implementing a machine learning process: the **how**

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations is achieved

Data distribution and I/O operations



Implementing a machine learning process: the **how**

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations is achieved

Data distribution and I/O operations



Implementing a machine learning process: the **how**

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations

Data distribution and I/O operations



Apache MESOS



APACHE Spark

APACHE Spark



Implementing a machine learning process: the **how**

The type of **statistical model**  MAHOUT

The **exact formula** being implemented  MAHOUT

The way the **algorithm** works  MAHOUT

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations

Data distribution and I/O operations



Apache MESOS



APACHE Spark

APACHE Spark



Implementing a machine learning process: the **how**

The type of **statistical model**  MAHOUT

The **exact formula** being implemented  MAHOUT

The way the **algorithm** works  MAHOUT

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations

Data distribution and I/O operations



Apache MESOS™



Implementing a machine learning process: the **how**

The type of **statistical model**



This is the Spark Machine Learning Library

The **exact formula** being implemented



The way the **algorithm** works



The way **individual calculations** are computed

The way the **distribution of the operations** is achieved

The way the **control of the distribution** of the operations



Apache MESOS

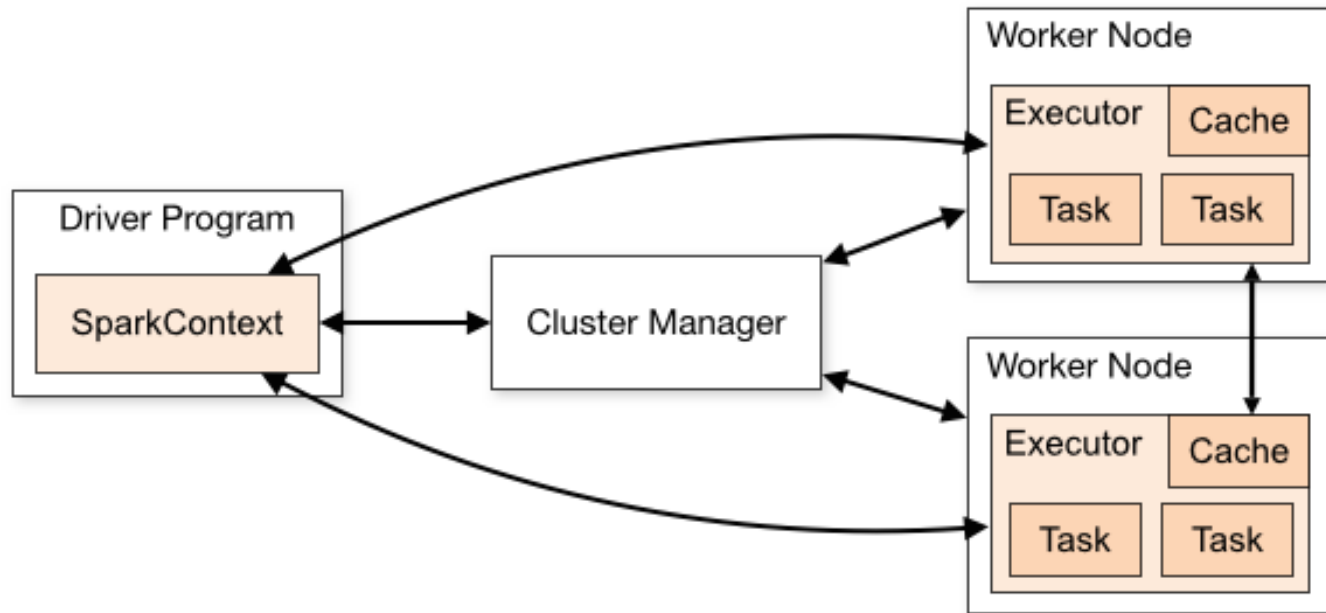


Data distribution and I/O operations

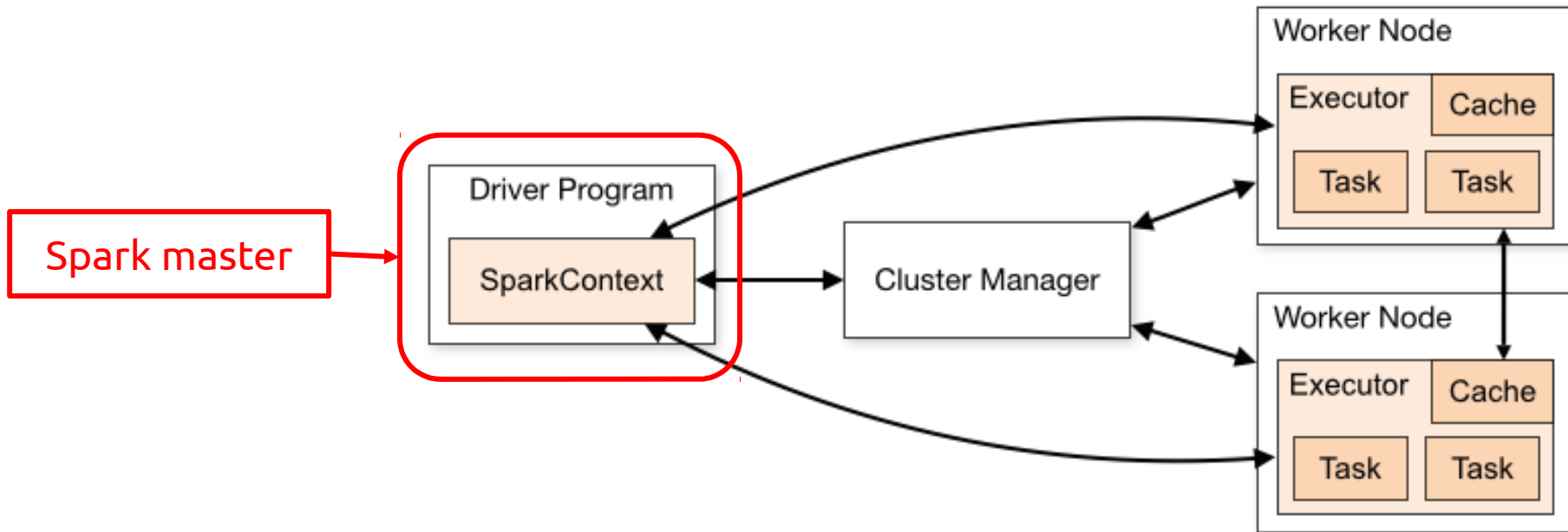


Apache Spark and Hadoop architectures in a nutshell

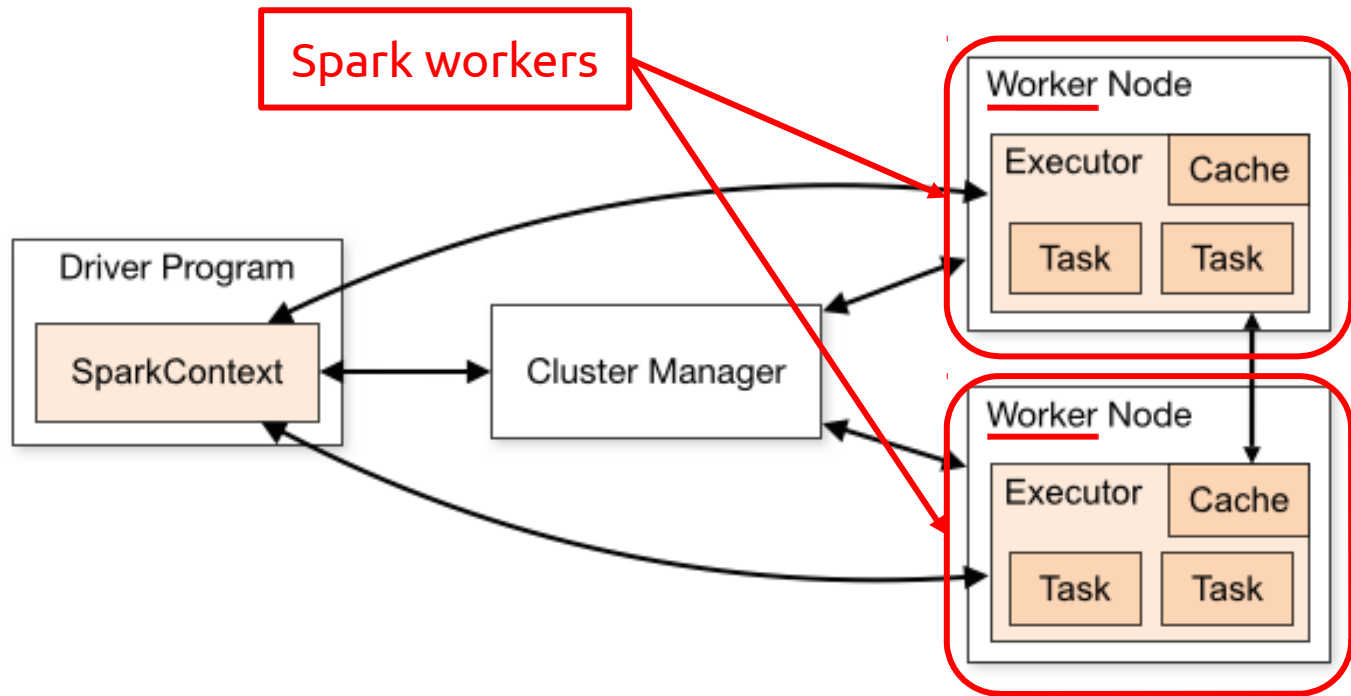
Spark architecture basics



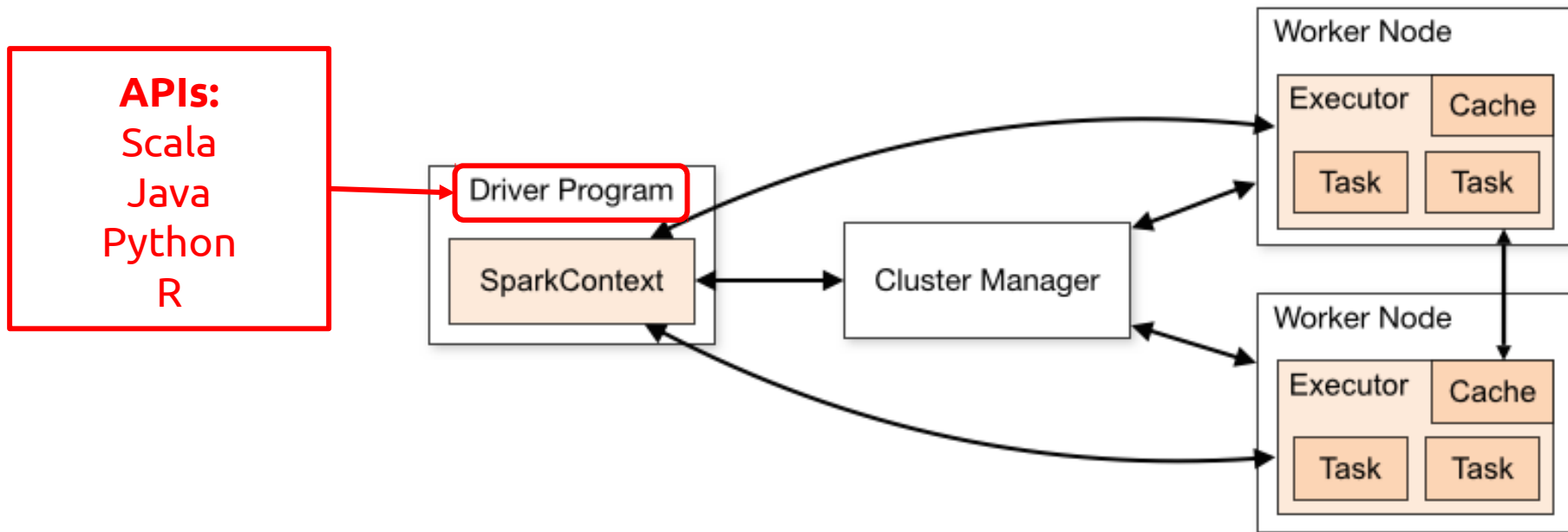
Spark architecture basics



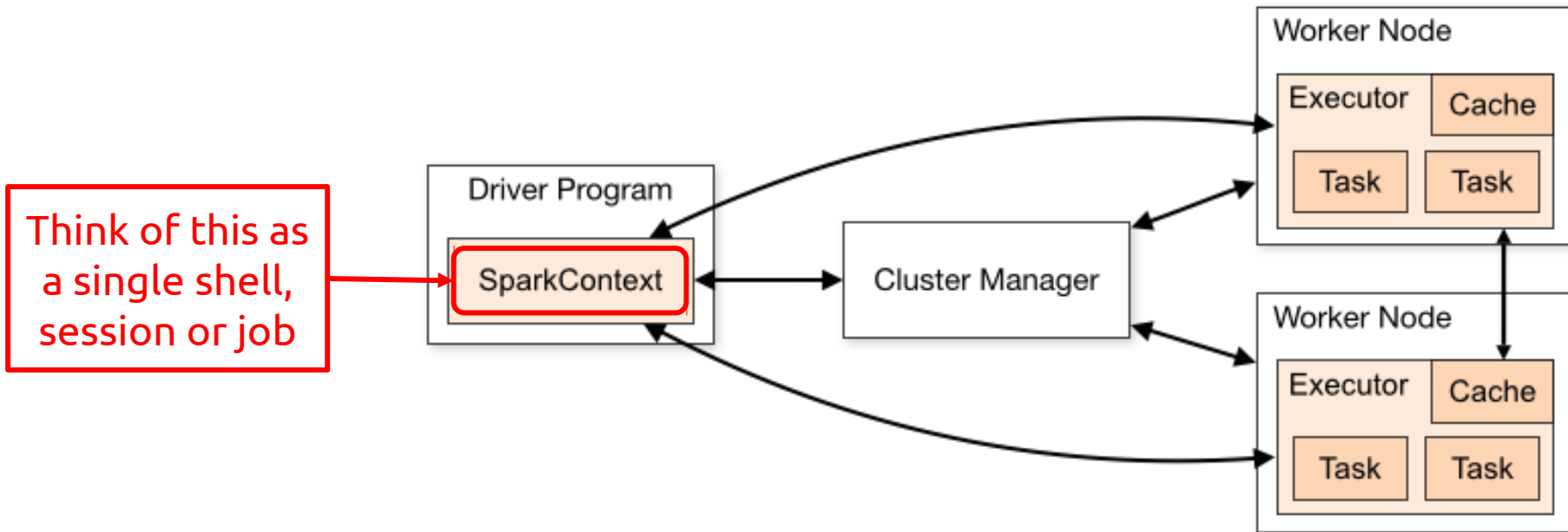
Spark architecture basics



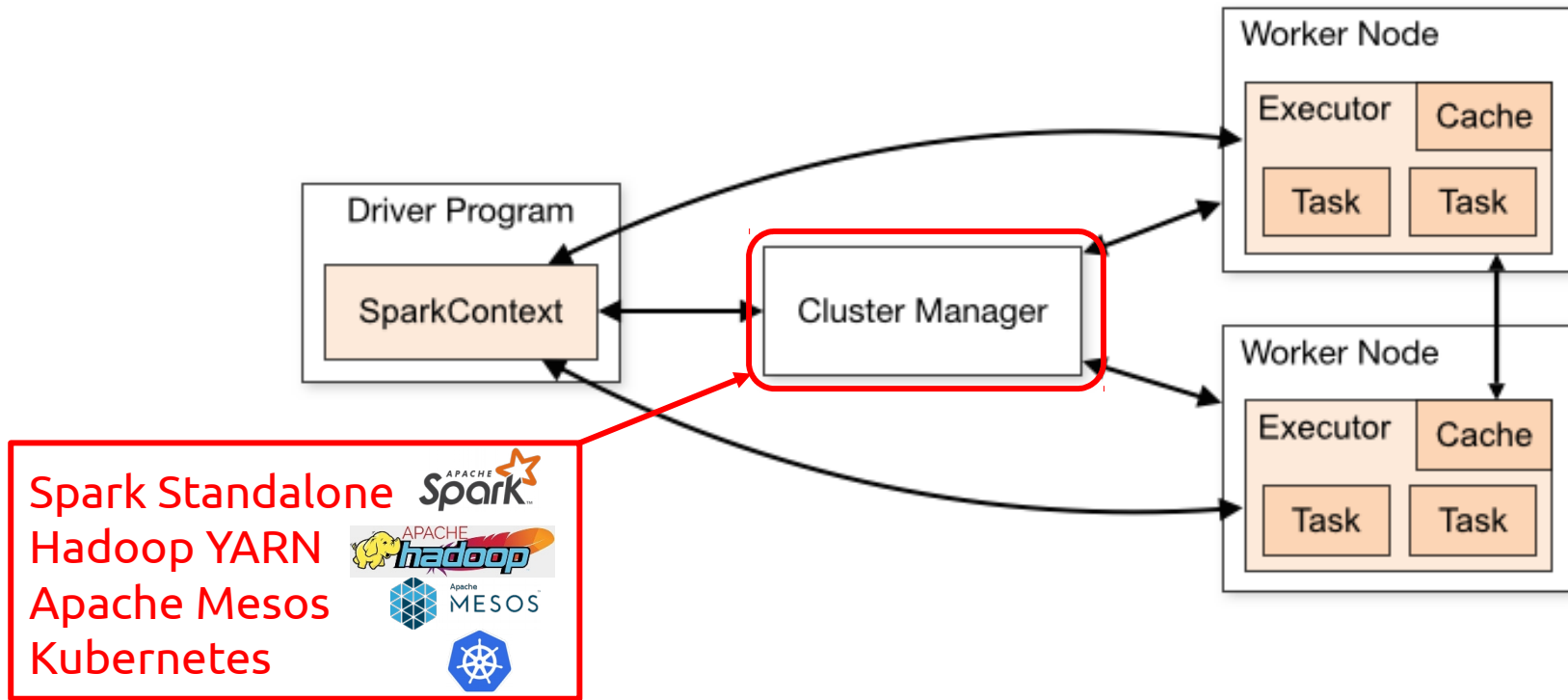
Spark architecture basics



Spark architecture basics



Spark architecture basics



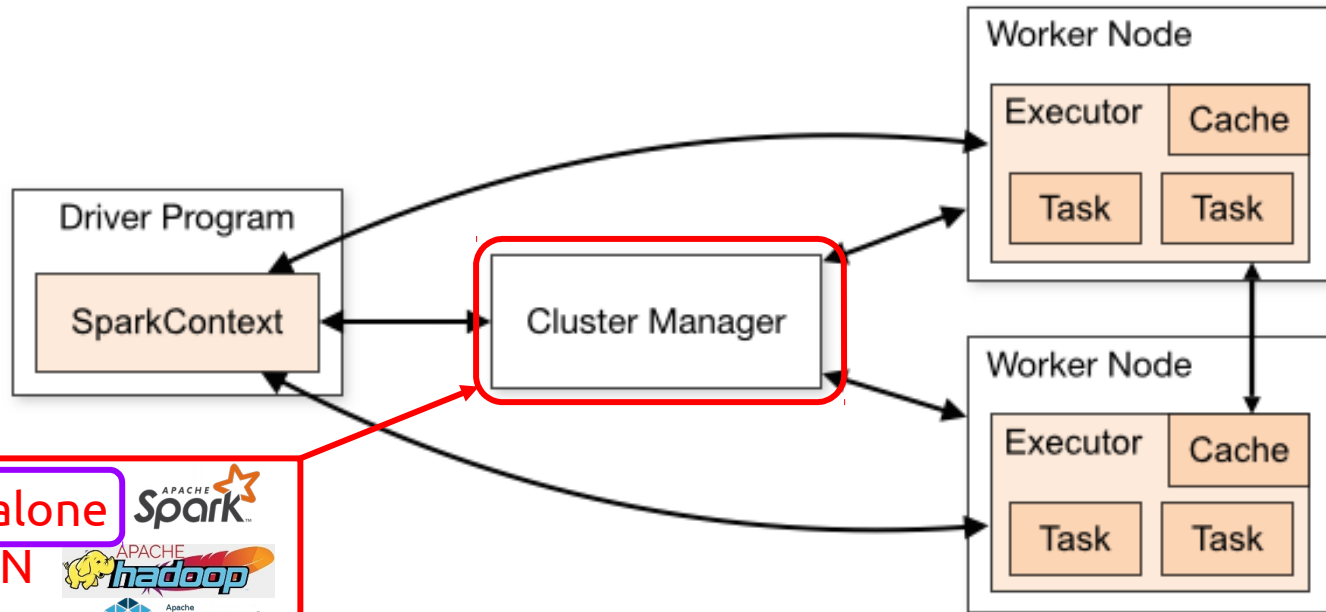
Spark Standalone
 Hadoop YARN
 Apache Mesos
 Kubernetes



Spark architecture basics

IMPORTANT:
Standalone
does not equal
local mode!

Spark Standalone 
 Hadoop YARN 
 Apache Mesos 
 Kubernetes 



Understanding Hadoop for Spark

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved: **Hadoop MapReduce**

The way the **control of the distribution** of the operations: **Hadoop YARN**

Data distribution and I/O operations: **Hadoop Distributed File System (HDFS)**

Understanding Hadoop for Spark

The type of **statistical model**

The **exact formula** being implemented

The way the **algorithm** works

The way **individual calculations** are computed

The way the **distribution of the operations** is achieved: **Hadoop MapReduce**

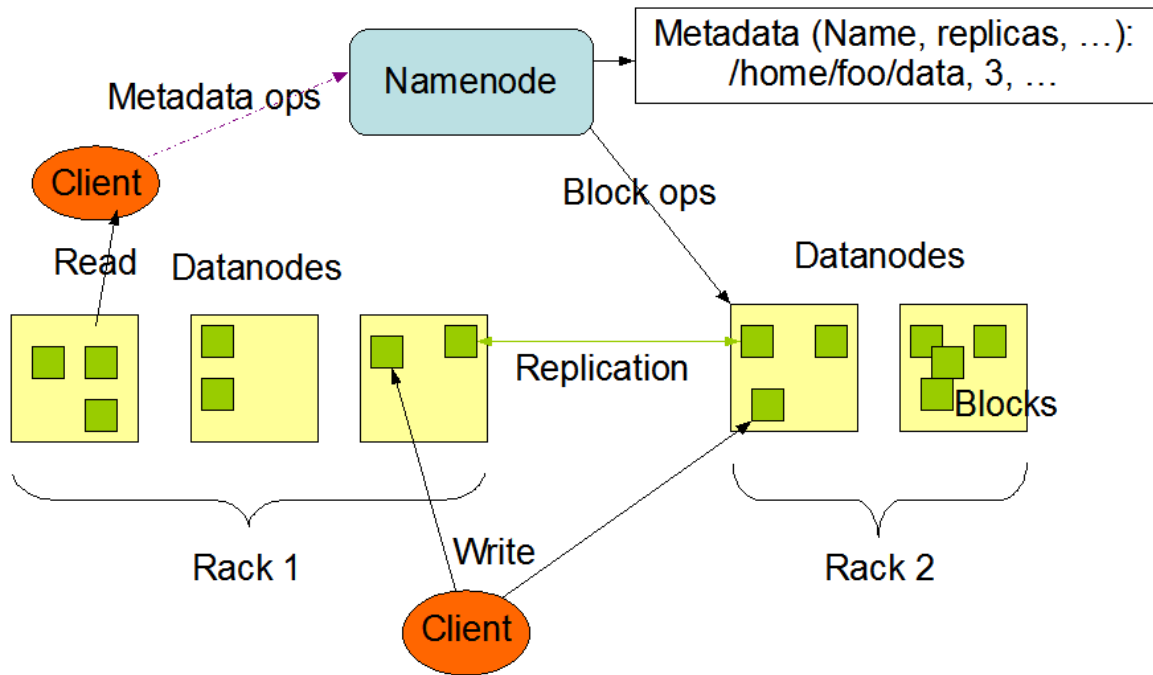
The way the **control of the distribution** of the operations: **Hadoop YARN**

Data distribution and I/O operations: **Hadoop Distributed File System (HDFS)**



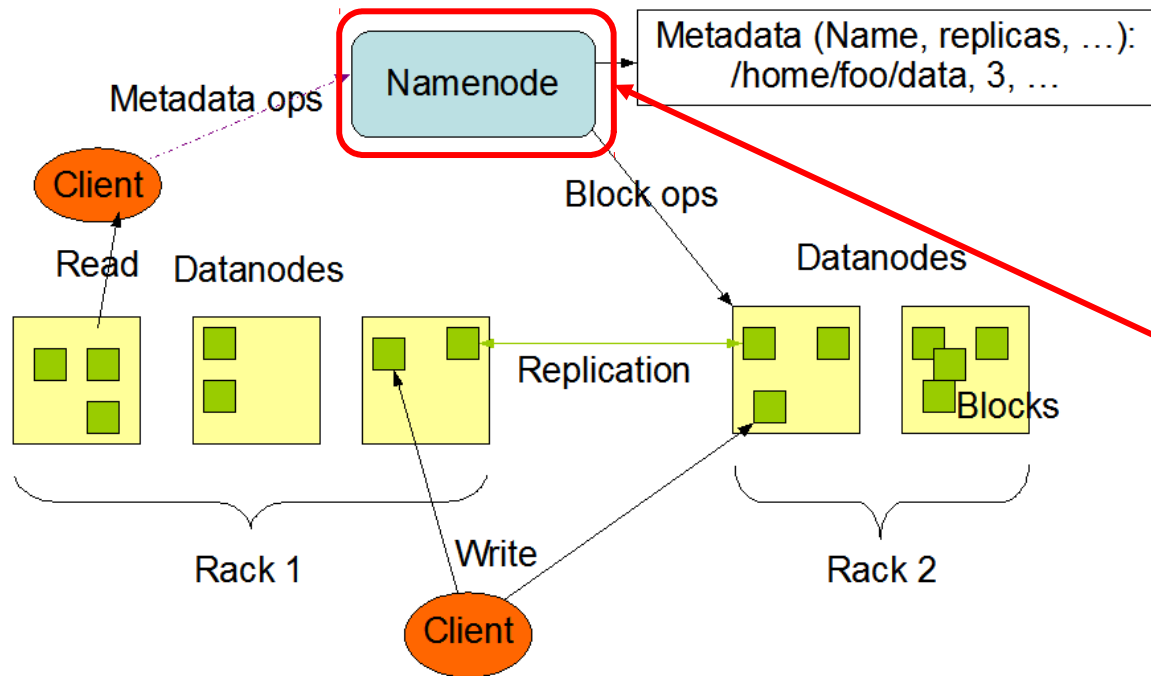
The Hadoop Distributed File System (HDFS)

HDFS Architecture



The HDFS NameNode

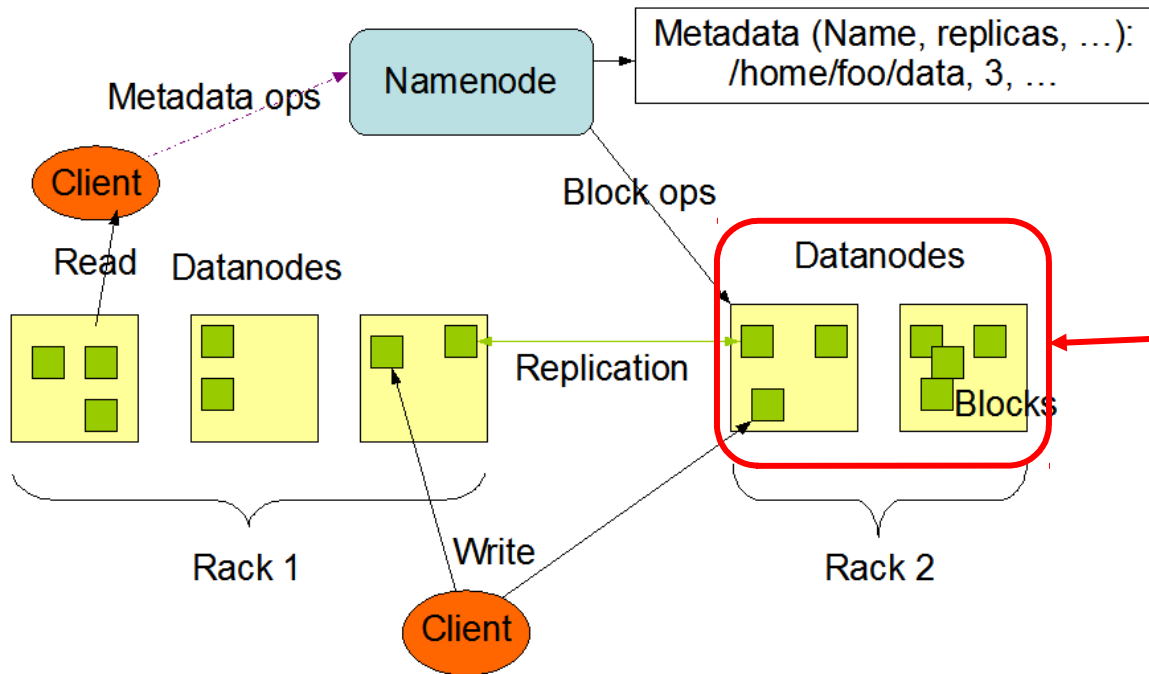
HDFS Architecture



- This is running **on the Spark master.**
- This is what we see when we browse the filesystem.

The HDFS DataNodes

HDFS Architecture



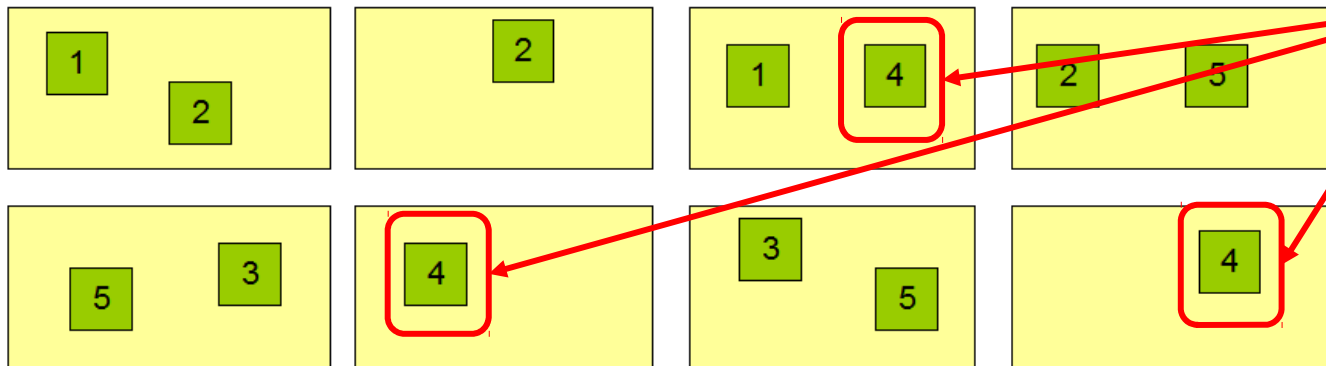
- These are running **on the Spark workers.**
- This is where the data actually is.

Data replication and fault tolerance in the HDFS

Block Replication

Namenode (Filename, numReplicas, block-ids, ...)
 /users/sameerp/data/part-0, r:2, {1,3}, ...
 /users/sameerp/data/part-1, r:3, {2,4,5}, ...

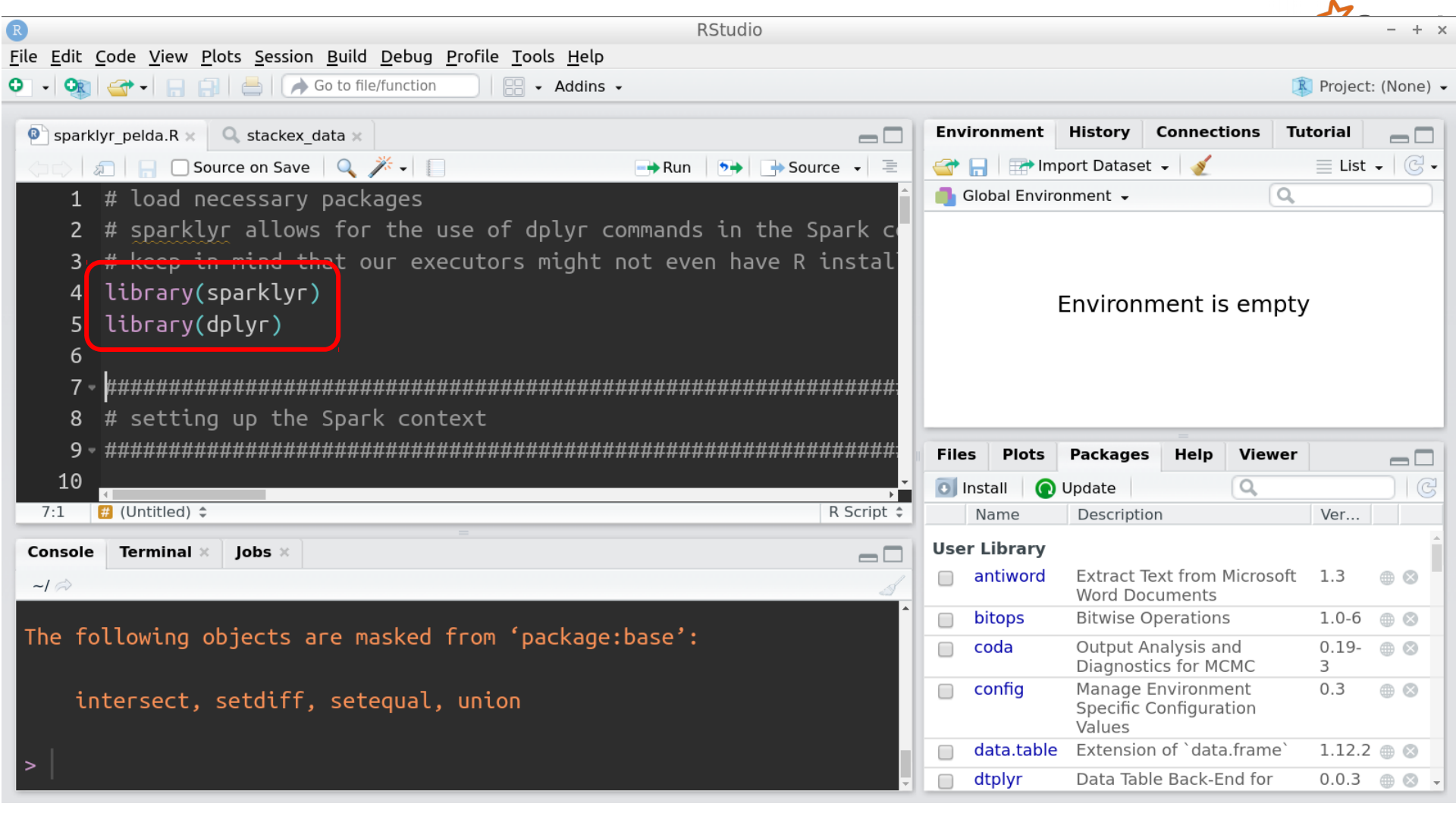
Datanodes



Data block replicated on different nodes

Let's get practical!

You can find all code and explanations at:
<https://github.com/zkpti/poltext2019-sparktutorial>



sparklyr_pelda.R x stackex_data x

Source on Save Run Source

```
1 # load necessary packages
2 # sparklyr allows for the use of dplyr commands in the Spark context
3 # keep in mind that our executors might not even have R installed
4 library(sparklyr)
5 library(dplyr)
6
7 #####
8 # setting up the Spark context
9 #####
10
```

7:1 # (Untitled) R Script

Console Terminal x Jobs x

~/

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

Environment History Connections Tutorial

Import Dataset List

Global Environment

Environment is empty

Files Plots Packages Help Viewer

Install Update

Name Description Ver...

User Library

<input type="checkbox"/>	antiword	Extract Text from Microsoft Word Documents	1.3	⊕ ⊗
<input type="checkbox"/>	bitops	Bitwise Operations	1.0-6	⊕ ⊗
<input type="checkbox"/>	coda	Output Analysis and Diagnostics for MCMC	0.19-3	⊕ ⊗
<input type="checkbox"/>	config	Manage Environment Specific Configuration Values	0.3	⊕ ⊗
<input type="checkbox"/>	data.table	Extension of `data.frame`	1.12.2	⊕ ⊗
<input type="checkbox"/>	dtplyr	Data Table Back-End for	0.0.3	⊕ ⊗

sparklyr_pelda.R x stackex_data x

Source on Save

Run

Source

```
15
16 # create Spark config
17 # this config list allows us to append further configuration of
18 conf <- spark_config()
19
20 # populate Spark config with further settings
21 # remember these numbers can and should change based on your m
22 # the amount of RAM to allocate to the Spark context, remember
23 # as a rule of thumb you can allow the Spark context to have s
24 # in the present case everything is on one local machine, so w
25
```

16:1 # (Untitled)

R Script

Console Terminal x Jobs x

```
~/
> # set system environment variables
> Sys.setenv(SPARK_HOME = "/opt/spark")
> # only need to set JAVA_HOME if you have more than one Java versio
n on your system
> Sys.setenv(JAVA_HOME = "/usr/lib/jvm/java-8-openjdk-amd64")
>
```

Environment History Connections Tutorial

Import Dataset













List

Global Environment

Environment is empty

Files Plots Packages Help Viewer

Install Update

	Name	Description	Ver...	
<input type="checkbox"/>	antiword	Extract Text from Microsoft Word Documents	1.3	 
<input type="checkbox"/>	bitops	Bitwise Operations	1.0-6	 
<input type="checkbox"/>	coda	Output Analysis and Diagnostics for MCMC	0.19-3	 
<input type="checkbox"/>	config	Manage Environment Specific Configuration Values	0.3	 
<input type="checkbox"/>	data.table	Extension of `data.frame`	1.12.2	 
<input type="checkbox"/>	dtplyr	Data Table Back-End for	0.0.3	 

sparklyr_pelda.R x stackex_data x

Source on Save

Run

Source

```
24 # in the present case everything is on one local machine, so w
25 conf$spark.executor.memory <- "2GB"
26 conf$spark.driver.memory <- "2GB"
27 # same applies to cores/cpus/vcpus, leave some for the OS and
28 conf$spark.executor.cores <- 4
29 conf$spark.driver.cores <- 1
30 # this is a setting that is required by some operations to avo
31 conf$spark.serializer <- "org.apache.spark.serializer.KryoSeri
32 conf$spark.kryoserializer.buffer <- "256m"
33 conf$spark.kryoserializer.buffer.max <- "256m"
```

30:1 # (Untitled)

R Script

Console Terminal x Jobs x

```
~/ ↵
> conf$spark.driver.memory <- "2GB"
> # same applies to cores/cpus/vcpus, leave some for the OS and othe
r programs too
> conf$spark.executor.cores <- 4
> conf$spark.driver.cores <- 1
>
```

Environment History Connections Tutorial

Import Dataset

Global Environment

Data

conf List of 9

Files Plots Packages Help Viewer

Install Update

	Name	Description	Ver...		
<input type="checkbox"/>	antiword	Extract Text from Microsoft Word Documents	1.3	⊕	⊗
<input type="checkbox"/>	bitops	Bitwise Operations	1.0-6	⊕	⊗
<input type="checkbox"/>	coda	Output Analysis and Diagnostics for MCMC	0.19-3	⊕	⊗
<input type="checkbox"/>	config	Manage Environment Specific Configuration Values	0.3	⊕	⊗
<input type="checkbox"/>	data.table	Extension of `data.frame`	1.12.2	⊕	⊗
<input type="checkbox"/>	dtplyr	Data Table Back-End for	0.0.3	⊕	⊗

```

40
47 # create the Spark context
48 # we have to define the master, which points to the Spark clus
49 # remember we know the address of our master from the initial
50 # setting the app_name is not necessary, but helps when later
51 # we pass our list of configurations to the Spark context crea
52 sc <- spark_connect(master="spark://master-neve:7077",
53                     app_name = "sparklyr-test-poltext",
54                     config = conf)
55
56
47:1 # (Untitled) R Script

```

```

> conf$spark.dynamicAllocation.enabled <- "false"
> # when starting the Spark context in local mode from an RStudio Se
rver (as opposed to a simple RStudio) the time to establish the conn
ection takes too long
> conf$spark.gateway.start.timeout <- 120
>

```

Data

conf	List of 19
------	------------

User Library

<input type="checkbox"/>	antiword	Extract Text from Microsoft Word Documents	1.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	bitops	Bitwise Operations	1.0-6	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	coda	Output Analysis and Diagnostics for MCMC	0.19-3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	config	Manage Environment Specific Configuration Values	0.3	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	data.table	Extension of `data.frame`	1.12.2	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	dtplyr	Data Table Back-End for	0.0.3	<input type="checkbox"/>	<input type="checkbox"/>


```

sparklyr_pelda.R x stackex_data x
Source on Save Run Source
50 # setting the app_name is not necessary, but helps when later
51 # we pass our list of configurations to the Spark context crea
52 sc <- spark_connect(master="spark://bland-ThinkPad-T490:7077",
53                      app_name = "sparklyr-test-poltext",
54                      config = conf)
55
56 # let's see if the Spark context is indeed running
57 # we can also open the Spark web UI from inside RStudio
58 spark_web(sc)
59
60
56:1 # (Untitled) R Script

```

```

Console Terminal x Jobs x
~/
we can start the spark context without this as well, in which case
it just uses the default values
> sc <- spark_connect(master="spark://bland-ThinkPad-T490:7077",
+                      app_name = "sparklyr-test-poltext",
+                      config = conf)
>

```

Environment History Connections Tutorial

spark://bland-ThinkPad-T490:7077 master-neve 0:7077 Spark

(No tables)

Files Plots Packages Help Viewer

Install Update

Name	Description	Versi...
<input type="checkbox"/> antiword	Extract Text from Microsoft Word Documents	1.3
<input type="checkbox"/> bitops	Bitwise Operations	1.0-6
<input type="checkbox"/> coda	Output Analysis and Diagnostics for MCMC	0.19-3
<input type="checkbox"/> config	Manage Environment Specific Configuration Values	0.3
<input type="checkbox"/> data.table	Extension of `data.frame`	1.12.2
<input type="checkbox"/> dtplyr	Data Table Back-End for	0.0.3

User Library

sparklyr_pelda.R x stackex_data x

Source on Save

Run

Source

```
63
64 # loading data into the Spark context from the HDFS is the best
65 stackex_data <- spark_read_csv(sc,
66                               "stackex_data",
67                               "~/MTA_PTI/ELKH_kepzes_2021_maj",
68                               header = TRUE,
69                               infer_schema = FALSE,
70                               delimiter = ";",
71                               memory = FALSE) # the default is true
72
```

73:1 # (Untitled)

R Script

Console Terminal x Jobs x

```
~/
+                               delimiter = ";",
+                               memory = FALSE) # the default is true,
ue, meaning the table will be automatically cached, but the table mi
ght be too big, or we might just need a few columns, so let's set th
is to false and cache the table ourselves, if we need to
>
```

Environment History Connections Tutorial

spark://l master-neve 0:7077 Spark

stackex_data

Files Plots Packages Help Viewer

Install Update

Name Description Versi...

User Library

<input type="checkbox"/>	antiword	Extract Text from Microsoft Word Documents	1.3	⊙ ⊗
<input type="checkbox"/>	bitops	Bitwise Operations	1.0-6	⊙ ⊗
<input type="checkbox"/>	coda	Output Analysis and Diagnostics for MCMC	0.19-3	⊙ ⊗
<input type="checkbox"/>	config	Manage Environment Specific Configuration Values	0.3	⊙ ⊗
<input type="checkbox"/>	data.table	Extension of `data.frame`	1.12.2	⊙ ⊗
<input type="checkbox"/>	dtplyr	Data Table Back-End for	0.0.3	⊙ ⊗

```
sparklyr_pelda.R x stackex_data x
63
64 # loading data into the Spark context from the HDFS is the best
65 stackex_data <- spark_read_csv(sc,
66                               "stackex_data",
67                               "~/MTA_PTI/ELKH_kepzes_2021_maj",
68                               header = TRUE,
69                               infer_schema = FALSE,
70                               delimiter = ";",
71                               memory = FALSE) # the default is
72
73:1 # (Untitled) R Script
```

```
Console Terminal x Jobs x
~/
+                               delimiter = ";",
+                               memory = FALSE) # the default is tr
ue, meaning the table will be automatically cached, but the table mi
ght be too big, or we might just need a few columns, so let's set th
is to false and cache the table ourselves, if we need to
>
```

Environment History Connections Tutorial

Import Dataset

Global Environment

Data

conf	List of 19	🔍
sc	List of 13	🔍
stackex_...	List of 2	🔍

Files Plots Packages Help Viewer

Install Update

Name	Description	Versi...
<input type="checkbox"/> antiword	Extract Text from Microsoft Word Documents	1.3
<input type="checkbox"/> bitops	Bitwise Operations	1.0-6
<input type="checkbox"/> coda	Output Analysis and Diagnostics for MCMC	0.19-3
<input type="checkbox"/> config	Manage Environment Specific Configuration Values	0.3
<input type="checkbox"/> data.table	Extension of `data.frame`	1.12.2
<input type="checkbox"/> dtplyr	Data Table Back-End for	0.0.3

sparklyr_pelda.R x stackex_data x

Show Attributes

Name	Type	Value
stackex_data	NULL	Pairlist of length 0

(No selection)

Console Terminal x Jobs x

```
~/   
+ memory = FALSE) # the default is tr   
ue, meaning the table will be automatically cached, but the table mi   
ght be too big, or we might just need a few columns, so let's set th   
is to false and cache the table ourselves, if we need to   
> View(stackex_data)   
> |
```

Environment History Connections Tutorial

Import Dataset

Global Environment

Data

▶ conf	List of 19	🔍
▶ sc	List of 13	🔍
▶ stackex_...	List of 2	🔍

Files Plots Packages Help Viewer

Install Update

Name Description Versi...

User Library

<input type="checkbox"/>	antiword	Extract Text from Microsoft Word Documents	1.3	⊕ ⊗
<input type="checkbox"/>	bitops	Bitwise Operations	1.0-6	⊕ ⊗
<input type="checkbox"/>	coda	Output Analysis and Diagnostics for MCMC	0.19-3	⊕ ⊗
<input type="checkbox"/>	config	Manage Environment Specific Configuration Values	0.3	⊕ ⊗
<input type="checkbox"/>	data.table	Extension of `data.frame`	1.12.2	⊕ ⊗
<input type="checkbox"/>	dtplyr	Data Table Back-End for	0.0.3	⊕ ⊗

sparklyr_pelda.R x stackex_data x

Show Attributes

Name	Type	Value
stackex_data	NULL	Pairlist of length 0

stackex_data

Environment History Connections Tutorial

Import Dataset

Global Environment

stackex_... List of 2

- src: List of 1
- ..\$ con: List of 13
 -\$ master : chr "spark://bland..."
 -\$ method : chr "shell"
 -\$ app_name : chr "sparklyr-te..."

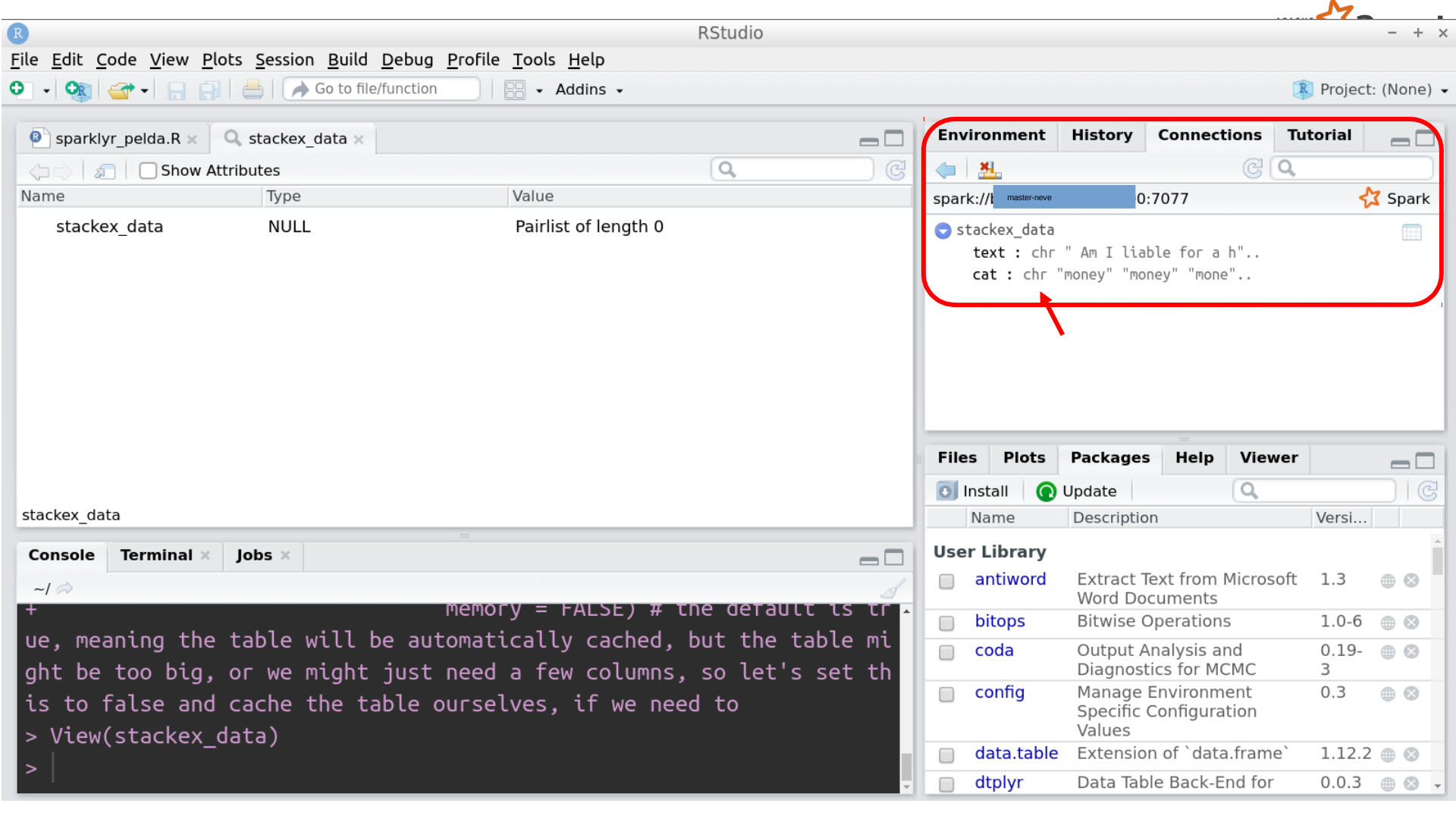
Console Terminal x Jobs x

```
~/  
+ ... memory = FALSE) # the default is tr  
ue, meaning the table will be automatically cached, but the table mi  
ght be too big, or we might just need a few columns, so let's set th  
is to false and cache the table ourselves, if we need to  
> View(stackex_data)  
>
```

Files Plots Packages Help Viewer

Install Update

Name	Description	Versi...
<input type="checkbox"/> antiword	Extract Text from Microsoft Word Documents	1.3
<input type="checkbox"/> bitops	Bitwise Operations	1.0-6
<input type="checkbox"/> coda	Output Analysis and Diagnostics for MCMC	0.19-3
<input type="checkbox"/> config	Manage Environment Specific Configuration Values	0.3
<input type="checkbox"/> data.table	Extension of `data.frame`	1.12.2
<input type="checkbox"/> dtplyr	Data Table Back-End for	0.0.3



```
sparklyr_pelda.R x stackex_data x
Source on Save Run Source
75
76 # you can also use the copy_to command (see https://cran.r-pro
77 # it is important to keep track of where our different data el
78
79 # check the data
80 sdf_nrow(stackex_data)
81 sdf_schema(stackex_data)
82 head(stackex_data,
83       n=2)
84
85
```

```
Console Terminal x Jobs x
~/
ght be too big, or we might just need a few columns, so let's set th
is to false and cache the table ourselves, if we need to
> # check the data
> sdf_nrow(stackex_data)
[1] 754
>
```

Environment History Connections Tutorial

spark://l master-neve 0:7077 Spark

stackex_data

Files Plots Packages Help Viewer

Install Update

Name	Description	Version
<input type="checkbox"/> antiword	Extract Text from Microsoft Word Documents	1.3
<input type="checkbox"/> bitops	Bitwise Operations	1.0-6
<input type="checkbox"/> coda	Output Analysis and Diagnostics for MCMC	0.19-3
<input type="checkbox"/> config	Manage Environment Specific Configuration Values	0.3
<input type="checkbox"/> data.table	Extension of `data.frame`	1.12.2
<input type="checkbox"/> dtplyr	Data Table Back-End for	0.0.3

```
sparklyr_pelda.R x stackex_data x  
79 # check the data  
80 sdf_nrow(stackex_data)  
81 sdf_schema(stackex_data)  
82 head(stackex_data,  
83       n=2)  
84  
85 stackex_data %>%  
86   group_by(cat) %>%  
87   summarise(n = n())  
88  
89
```

Console Terminal x Jobs x
~/
+ n=2)
Source: spark<?> [?? x 2]
text cat
<chr> <chr>
1 " Am I liable for a health insurance bill reopened I ... money
2 " Assuming I know X questions We ve had a few questi... money

Environment History Connections Tutorial

spark:// master-neve 0:7077 Spark

stackex_data
text : chr " Am I liable for a h..
cat : chr "money" "money" "mone"..

Files Plots Packages Help Viewer

Install Update

Name	Description	Versi...
<input type="checkbox"/> antiword	Extract Text from Microsoft Word Documents	1.3
<input type="checkbox"/> bitops	Bitwise Operations	1.0-6
<input type="checkbox"/> coda	Output Analysis and Diagnostics for MCMC	0.19-3
<input type="checkbox"/> config	Manage Environment Specific Configuration Values	0.3
<input type="checkbox"/> data.table	Extension of `data.frame`	1.12.2
<input type="checkbox"/> dtplyr	Data Table Back-End for	0.0.3


```
sparklyr_pelda.R x stackex_data x
Source on Save Run Source
83     n=2)
84
85     stackex_data %>%
86       group_by(cat) %>%
87       summarise(n = n())
88
89 # setting up the text preprocessing pipeline
90 # NOTE: all the ft_ functions correspond to and are wrappers for
91 preproc_pipeline <- ml_pipeline(sc) %>%
92   # tokenize on white space and transform to lowercase, also d
93
```

89:1 # (Untitled)

R Script

Console Terminal x Jobs x

~/

```
cat          n
<chr>       <dbl>
1 anime      103
2 android    163
3 money       48
4 retrocomputing 187
```

Environment History Connections Tutorial

spark://k master-neve 0:7077 Spark

stackex_data

Files Plots Packages Help Viewer

Install Update

Name	Description	Versi...
<input type="checkbox"/> antiword	Extract Text from Microsoft Word Documents	1.3
<input type="checkbox"/> bitops	Bitwise Operations	1.0-6
<input type="checkbox"/> coda	Output Analysis and Diagnostics for MCMC	0.19-3
<input type="checkbox"/> config	Manage Environment Specific Configuration Values	0.3
<input type="checkbox"/> data.table	Extension of `data.frame`	1.12.2
<input type="checkbox"/> dtplyr	Data Table Back-End for	0.0.3



Spark Jobs (?)

User: [redacted]
Total Uptime: 4.9 min
Scheduling Mode: FIFO
Active Jobs: 1
Completed Jobs: 11, only showing 10

▶ Event Timeline

▼ Active Jobs (1)

Job Id ▼	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
10	collect at utils.scala:204 collect at utils.scala:204 (kill)	2021/05/24 16:52:29	5 s	1/2	70/78 (5 running)

▼ Completed Jobs (11, only showing 10)

Job Id ▼	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
9	collect at utils.scala:204 collect at utils.scala:204	2021/05/24 16:52:21	7 s	1/1 (1 skipped)	100/100 (4 failed) (3 skipped)

```
91 preproc_pipeline <- ml_pipeline(sc) %>%
92   # tokenize on white space and transform to lowercase, also drop single characters
93   ft_regex_tokenizer(input_col = "text",
94                       output_col = "words",
95                       min_token_length = 2,
96                       to_lower_case = TRUE) %>%
97   # remove stopwords
98   ft_stop_words_remover(input_col = "words",
99                         output_col = "words2",
100                        stop_words = ml_default_stop_words(sc,
101                                                            language = "english")) %>%
102   # create term frequency vector (drop word if document frequency is less than 5)
103   ft_count_vectorizer(input_col = "words2",
104                       output_col = "raw_features",
105                       min_df = 5) %>%
106   # create term frequency weighted by inverse document frequency vector
107   ft_idf(input_col = "raw_features",
```

Environment
is empty

<input type="checkbox"/>	Ex	⊗
<input type="checkbox"/>	Te:	
<input type="checkbox"/>	frc	
<input type="checkbox"/>	Mi	
<input type="checkbox"/>	Wc	
<input type="checkbox"/>	Dc	
<input type="checkbox"/>	Bit	⊗
<input type="checkbox"/>	Oç	
<input type="checkbox"/>	Ou	⊗
<input type="checkbox"/>	An	

```
104         output_col = "raw_features",
105         min_df = 5) %>%
106     # create term frequency weighted by inverse document frequency vector
107     ft_idf(input_col = "raw_features",
108           output_col = "features")
109
110 # call pipeline to transform the data
111 stackex_data <- ml_fit_and_transform(preproc_pipeline,
112                                     stackex_data)
113
114 # check the data again
115 sdf_nrow(stackex_data)
116 sdf_schema(stackex_data)
117 head(stackex_data,
118       n=2)
119
```

```
121 # NOTE: all the ml_ functions correspond to and are wrappers for Spark Mlib functions, ju
122 nb_pipeline <- ml_pipeline(sc) %>%
123   # have to change the text labels to numeric for the classifier to be able to handle it
124   ft_string_indexer(input_col = "cat",
125                     output_col = "label") %>%
126   # add the naive bayes classifier to the pipeline
127   ml_naive_bayes()
128
129 # parameter grid for parameter tuning (this is very small and simple on purpose for the tu
130 # the name of the list element that corresponds to a pipeline element (in this case "naive
131 # you can add more than one parameter to tune at the same time, note how the parameter nam
132 param_grid <- list(
133   naive_bayes = list(
134     smoothing = c(1.0,
135                  0.5)
136   )
137 )
```

```
139 # set up cross validator for parameter tuning
140 # note: the default metric for ml_multiclass_classification_evaluator is f1, we could also
141 # see https://cran.r-project.org/web/packages/sparklyr/sparklyr.pdf
142 cv <- ml_cross_validator(sc,
143     estimator = nb_pipeline,
144     estimator_param_maps = param_grid,
145     evaluator = ml_multiclass_classification_evaluator(sc),
146     num_folds = 3,
147     parallelism = 4 # this is based on the number of cores/cpus/vcpus
148 )
149
150 # create train-test split
151 split_data <- sdf_random_split(stackex_data,
152     training = 0.7,
153     test = 0.3)
```



```
193 # since this can be too much to handle, let's just collect the columns we need for now
194 test_with_pred_redux <- select(test_with_pred,
195                               cat,
196                               pred_cat)
197 test_df <- sdf_collect(test_with_pred_redux)
198
199 install.packages("tidyr") # we need this for the spread function
200 library(tidyr)
201
202 # now we can use the spread function on our dataframe in the R session to create the con
203 test_df %>%
204   group_by(cat,
205            pred_cat) %>%
206   summarise(n = n()) %>%
207   spread(key = cat,
208          value = n)
```



```

214 # (just as it is important to keep track of where our data is, it is also important to keep
215 test_with_pred %>%
216   group_by(cat,
217             pred_cat) %>% # this happens in the Spark context
218   summarise(n = n()) %>% # this happens in the Spark context
219   collect() %>% # this is where we move our data from the Spark context to the R
220   spread(key = cat,
221           value = n) # this happens in the R session
222
223 # hmmm, maybe we should check our metrics without the gardening category
224 test_with_pred %>%
225   filter(label != 0) %>%
226   ml_multiclass_classification_evaluator(label_col = "label",
227                                           prediction_col = "prediction",
228                                           metric_name = "f1")
229
230 # let's try the linear SVC model

```

Thank you for your attention!

You can find all code and explanations at:
<https://github.com/zkpti/poltext2019-sparktutorial>