

Data Engineering

Assignment 1: Big Data in Ihrem Umfeld

Schemalose Daten:

Ein Kunde für den ich arbeite verkauft Personendaten, diese sind in XML Dokumente gespeichert und können beliebig erweitert werden. Zum Beispiel durch neue Lifestyle Attribute usw. Als DB wird Marklogic verwendet.

Strukturierte Daten:

Parallel zur Marklogic DB gibt es noch eine Postgres Datenbank in der alle Kunden Informationen, Job Status usw gespeichert werden. Hier gibt es so gut wie keine Änderungen.

Batch Processing:

In den oben genannten Big Data Projekt werden die Daten durch Talend ETL Jobs aufbereitet und danach werden die Dokumente erstellt bzw aktualisiert. Der Verarbeitung passiert in Batches.

Streaming:

Am anderen Ende sitzt der Kunde der sich über eine Weboberfläche (Java + HTML + JS) diverse Selektionen zusammenbaut und diese gegen die Datenbank ausführt und sofort ein Ergebnis erhält.

Assignment 2: Big Data in Ihrem Umfeld

Ich habe mich für Apache Spark entschieden, da wir in der Firma mehr mit Spark arbeiten als mit Flink. Weiters ist in Spark Batch und Streaming möglich.

```
ceda@ceda-PC /cygdrive/c/Users/ceda/Downloads/spark-1.6.2-bin-hadoop2.6/bin
$ ./spark-shell.cmd
log4j:WARN No appenders could be found for logger (org.apache.hadoop.metrics2.lib.MutableMetricsFactory).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
Using Spark's repl log4j profile: org/apache/spark/log4j-defaults-repl.properties
To adjust logging level use sc.setLogLevel("INFO")
Welcome to

  ____      __
 / ___/____/  /
/  /_  /_  /  /
/   /  /  /  /
/_/_/  /_/_/  /_/_/

version 1.6.2

Using Scala version 2.10.5 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_71)
Type in expressions to have them evaluated.
Type :help for more information.
Spark context available as sc.

16/06/30 21:23:15 WARN General: Plugin (Bundle) "org.datanucleus.store.rdbms" is already registered. Ensure you dont have multiple JAR versions of the same plugin in the classpath.
16/06/30 21:23:15 WARN General: Plugin (Bundle) "org.datanucleus.api.jdo" is already registered. Ensure you dont have multiple JAR versions of the same plugin in the classpath.
16/06/30 21:23:15 WARN General: Plugin (Bundle) "org.datanucleus" is already registered. Ensure you dont have multiple JAR versions of the same plugin in the classpath.
16/06/30 21:23:15 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
16/06/30 21:23:17 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
```

Da ich ursprünglich aus dem Java Bereich komme, ist in meinem Toolstack immer Eclipse dabei. Und um auch auf Windows alle Vorteile einer Command Line nutzen zu können ist das wichtigste Tool Cygwin, da ich mir dadurch einige Eclipse Plugins erspare (GIT, MAVEN, ...) und andere Tools z.B. Putty erspare.

Data Science

Assignment 1: Technologien

Alternativen zu R und Python:

- Java
- SQL
- SAS

Data Science Tools:

- SQL
- Java
- Excel
- Shell
- Notepad ++
- XML → XPath, XQuery

Assignment 2: Technologien

Ich habe mich für Python entschieden, da es eine größere Community gibt, das heißt besserer Support, mehr Libs usw. Weiters kann man Python gut mit Webapplikation und Datenbank verwenden. R ist aus meiner Sicht doch mehr für Standalone Anwendungen zu gebrauchen.

```
C:\Users\ChristophEder>python
Python 3.5.2 (v3.5.2:4def2a2901a5, Jun 25 2016, 22:01:18) [MSC v.1900 32 bit (Intel)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> print "Python is fun!"
File "<stdin>", line 1
    print "Python is fun!"
    ^
SyntaxError: Missing parentheses in call to 'print'
>>> print "Python is fun!"
File "<stdin>", line 1
    print "Python is fun!"
    ^
SyntaxError: Missing parentheses in call to 'print'
>>> print "Python is fun!"
File "<stdin>", line 1
    print "Python is fun!"
    ^
SyntaxError: Missing parentheses in call to 'print'
>>> print("Python is fun!")
Python is fun!
>>>
```

Toolchain:

Ich würde PyCharm verwenden, da ich gute Erfahrungen mit den JetBrains Produkten bis jetzt gemacht habe (PHPStorm, IntelliJ, AndroidStudio). Ich finde es sehr praktisch, dass man von Haus aus bei den JetBrains Tools die Eclipse Shortcuts einstellen kann. Somit hat man in jeder IDE die gleichen Shortcuts.