# Kristiania University College

# Course Assignment
## Master of
## Applied Computer Science
Department of Technology

| | |
|---|---|
| Assignment title | Exploring Bitcoin |
| Course code | MA120 |
| Course name | Big Data |
| Due date | 18 Oktober |

**Declaration:**

Through the submission of this assignment, we hereby declare that this report is the result of our own work, and that all sources have been properly cited to throughout the text.

| | |
|---|---|
| Names of students | Theodore Midtbø Alstad ; Howie Chen |
| Student ID numbers | 865317 ; 866354 |

# Exploring Bitcoin

## ABSTRACT

**In this assignment we parsed a dataset, in the form of XMLs, through a Hadoop- and Pig apache. We learned how to write MapReduce jobs and Pig scripts, and how effective it is on our dataset based on XML files from Bitcoin stack exchange.**

## Introduction

We choose to work together because both of us have python background, therefor we choose python as main programming language. We explored Bitcoin as given dataset from [archive.org/download/stackexchange/bitcoin.stackexchange.com.7z](archive.org/download/stackexchange/bitcoin.stackexchange.com.7z) through Apaches: Hadoop and Pig.

## Main functions

The datasets contains only XML files, which it means values to different attributes have a lot of ascii characters, punctiations, numbers and HTML tags. We also needed a function to parse the XML files to Hadoop and to create a mapper for the data. Therefor we created three different main functions thats being reused through the project. :

### CleanBody

In cleanBody function it formats strings for mapper function. The function makes non-case sensitivity, removes ascii characters , HTML formatting, and treats anything saperated by blanked space or as separated words. Most of the words after separating/splitting will be counted as separate words, and this will effect the results of some tasks. For instance the name "Jens-Petter" will be counted as two words; "Jens" and "Petter". The input for the function is a string body which will be formatted and it returns a list of formatted string.

```
1  def cleanBody(body):
2      body = body.lower()
3      body = ascii(body)
4      body = sub("<.+?>","",body)
5      body = body.replace("/", " ")
6      body = body.strip()
7
8      for i in ignore_char:
9          body = body.replace(i, "")
10
11     body = body.split(" ")
12
13     return body
```

**Listing 1.** CleanBody function

### Mapper Core

mapper_core is the core of the mapper function; it prints out the relevant data in a format parseable by Hadoop. It functions through three modes, as determined by the parameter mode: "single", "double", and "triple".

- Single: Assumes input "words" is a list of words to print. Prints word in words as (word, 1). Ignores empty strings and spaces.

- Double: Assumes input "words" is two nested lists to print. Prints word, count in words as (word, count). Ignores empty strings and spaces.

- Triple: Assumes input "words" is three nested lists to print. Prints id, score, title in words as (id, score, title). does not ignore empty strings and spaces.

1

```
1  def mapper_core(words, mode="single"):
2      if mode == "single":
3          for word in words:
4              if word not in ["", " "]:
5                  print("%s %s" %(word,1))  #Emit the word
6
7      elif mode == "double":
8          in1, in2 = words
9          for word, count in zip(in1,in2):
10             if word not in ["", " "]:
11                 print("%s %s" %(word,count))  #emit the words
12
13     elif mode == "triple":
14         in1, in2,in3  = words
15         for id, score, title in zip(in1,in2,in3):
16             print("%s %s %s" %(id,score, title))  #emit the words
```

**Listing 2.** mapper_core function

### Xmlparser

The chosen method to interpret the dataset is to parse to the mapper function. Although parsing an entire XML-file takes up significant memory, this method fits our dataset. It has been separated out as a function so it may be easily replaced by other methods more fit for large files. The input for this function is a XML-file and the output is a parsed XML-file

```
1  def xmlparser(infile):
2      if not isinstance(infile, str):
3          infile = infile.detach()
4      mytree = ET.parse(infile)
5      myroot = mytree.getroot()
6      return myroot
```

**Listing 3.** xmlparser function

## Task 1 Warmup

### 1a) WordCount
**Assumption**: Count the words in body of questions PostTypeID="1" in the *Posts.xml* file. The result should be how many times a word occur in the body of questions.

**Implementation**The will

**Notes/Reflection** bye bye

```
1  abovennthanks 1
2  abovennthanksn 2
3  abovennwhat 1
4  abovennwill 1
5  aboventhe 1
6  aboventrying 1
7  abr 1
8  abra 4
9  abras 1
10 abreast 1
11 abridge 1
```

1a_output

### 1b) Unique words
**Assumption**: Write a MapReduce job where the result should be unique words in the titles PostTypeID="1" in the *Posts.xml* File.

**Implementation**

**Notes/Reflection**

```
1  asm
2  asn
3  aspect
4  aspects
5  aspnet
6  assemble
7  assembled
8  assembles
9  assembling
10 assert
11 assertion
```

1b_output

## 1c) MoreThan10

**Assumption**: Write a simple python code to check the title length in *Posts.xml*. The result should output how many titles have more than 10 words.

**Implementation**

**Notes/Reflection**

```
1  7600
```

1c_output

## 1d) Stopwords

**Assumption**: Write a simple python code based on task 1a to exclude stopwords from body of questions PostTypeID="1" in the *Posts.xml*. The output should be text file without any stopwords.

**Implementation**

**Notes/Reflection**

```
1  ati 8
2  atiradeon 1
3  atlcoin 1
4  atm 18
5  atms 5
6  atom 1
7  atomic 17
8  atomically 1
9  attach 8
10 attached 6
11 attaching 2
```

1d_output

## 1e) Pig top 10

**Assumption**: Write a pig script to select top 10 listed words after removing the stopwords from *Posts.xml*. The output should print out top 10 listed words and the corresponding occurrence rate.

**Implementation**

**Notes/Reflection**

```
1  bitcoin 5918
2  transaction 2376
3  wallet 2320
4  address 1550
5  block 1222
6  mining 1172
7  blockchain 1169
8  the 1148
```

```
 9  transactions 1127
10  bitcoins 1050
```

<div align="center">1e_output</div>

### 1f) Tags
**Assumption**: Write a MapReduce job to create a dictionary over unique tags in *Posts.xml*. The result should print the unique tags.

**Implementation**

**Notes/Reflection**

```
 1  moneylaundering
 2  moneysupply
 3  moneytransfer
 4  movecmd
 5  msigna
 6  mtgox
 7  multibit
 8  multibithd
 9  multifactor
10  multigateway
11  multipartycomputation
```

<div align="center">1f_output</div>

# Task 2 Discover

### 2a) Counting
**Assumption**: We chose to write a MapReduce job to count the total unique users there are in *Users.xml*. The result will print out how many unique users there are.

**Implementation**

**Notes/Reflection**

```
 1  56451
```

<div align="center">2a_output</div>

### 2b) Unique users
**Assumption**: Write a MapReduce job based on task 2a) to create a mapper and a reducer functions in *Users.xml*. The result should contain unique users in the dataset.

**Implementation**

**Notes/Reflection**

```
 1  10529
 2  1053
 3  10530
 4  10531
 5  10533
 6  10534
 7  10535
 8  10536
 9  10537
10  10538
11  10539
```

<div align="center">2b_output</div>

### 2c) Top miners

**Assumption**:Write a MapReduce job to find top 10 users based on attribute Reputation="x" in *Users.xml*. The result will print out top 10 users based on their reputation.

### Implementation

### Notes/Reflection

```
1  david 87354
2  user 64135
3  pieter 51598
4  andrew 46149
5  murch 36524
6  thepiachu 31489
7  nick 29072
8  stephen 27509
9  chris 24736
10 nate 21620
```

2c_output

### 2d Top) questions

**Assumption**: Write a MapReduce job to find top 10 title questions PostTypeID="1" based on attribute Score="x" in *Posts.xml*.The result lists top 10 questions using id, question and the score.

### Implementation

### Notes/Reflection

### 2e) Favorite questions

**Assumption**: Write a MapReduce job to find top 10 title questions PostTypeID="1" based on attribute FavoriteCount="x" in *Posts.xml*.The result lists top 10 questions like id, question and the score.

### Implementation

### Notes/Reflection

### 2f) Average answers
**Assumption**:

### Implementation

### Notes/Reflection

### 2g) Countries
**Assumption**: We chose to write a MapReduce job to discover users by countries in *Users.xml*. The result lists different countries and corresponding users.

### Implementation

### Notes/Reflection

### 2h) Names
**Assumption**We chose to write a MapReduce job to find the most popular names in *Users.xml*. The result lists top 10 common names and how many.

**Implementation**

**Notes/Reflection**

**2i) Answers**
**Assumption**: Write a simple python code to find how many titles of questions PostTypeId="x" have at least one answers based on attribute AnswerCount in *Users.xml*. The result prints out how many questions have been answered.

**Implementation**

**Notes/Reflection**

# Task 3 Numbers

**3a) Bigram**
**Assumption**: We chose to write a MapReduce job to find the most common pair of adjacent words in *Posts.xml*.For instance, "big data" or "Fast car" are examples of bigram. The result print the most common bigram

**Implementation**

**Notes/Reflection**

**3b) Trigram**
**Assumption**: This task is based on 3a, to find three words that appear consecutively in *Posts.xml*. The result print the most common trigram.

**Implementation**

**Notes/Reflection**

**3c) Combiner**
**Assumption**: Write a MapReduce job and add a combiner before the data sent to the reducer. The result should be the same as a reducer, but the reducer recive a smaller volume of data.

**Implementation**

**Notes/Reflection** Having a combiner in A MapReduce jobs saves us bandwidth and computational strain by decreasing the volume of data sent from the mapper. What combiner means is to have a semi-reducer after the mapper before sending it to the reducer.

**3d) Useless**
**Assumption**: We chose write a MapReduce job to find how many times the word "useless" in *Posts.xml* occurs in the body of questions PostTypeId="1". The result print how many times the word useless occurs.

**Implementation**

**Notes/Reflection**

# Task 4 Search engine

**4a) Title index**
**Assumption**: We chose to write MapReduce job to create index over titles, bodies and answers of questions in *Posts.xml*. We are after having a simple index that lists publications in which a search term/s and occur/s. The result is a list of words and their index appearance in posts

**Implementation**

### Notes/Reflection

```
1  angel, ,11695,8519,5197
2  angeles, ,5197,11695
3  angellist, ,11695
4  angels, ,11695,687
5  anger, ,2749
6  angered, ,11912
7  angle, ,2937,4110,1224
8  angles, ,1859
9  angry, ,2240,10044
```

4a_output

## Conclusion