

---

Please document your code, without sufficient documentation you won't receive any points.

## 1 Cypher Query Language (60 P)

We have already imported a large Wikipedia dataset into Neo4j, you will all work on the same Neo4j instance on Abu2. Write queries for the tasks below and execute them using the web interface.

- Print the names of all articles that link to the article *Artificial\_intelligence*.
- Find all shortest paths that can lead you from the article *Alan\_Turing* to the article *Gene\_Amdahl* by following Wikipedia page-links.
- How many articles exist in the dataset?
- Which article is the most strongly connected – characterized by the sum of all links (in+out).

### 1.1 Hints

Write access is deactivated to Neo4j to avoid accidental deletions. You can access the Neo4j web interface on port 7475 on abu2. Forward the port to your local machine for easy access: `ssh -L 7475:abu2:7475 cluster`. Then use the address `localhost:7475` in your browser to access Neo4j.

To learn more about Neo4j, use the query:

```
:play start
```

To understand the existing data, use the Cypher query:

```
MATCH (n)-->(m) RETURN n,m limit 10;
```

### Submission:

1-queries.txt Your cypher queries.

## 2 Cypher Data Model (60 P)

Extend the already provided data model from *Task 1*: 1) Add the article text and date of last change. 2) Add Wikipedia authors who can each have authored multiple articles and have a user name and email-address.

Describe the datamodel and provide some queries to insert example entities. Create a screenshot of your example graph. You may use <http://console.neo4j.org/> to create an example graph.

### Submission:

2-datamodel.txt Your cypher queries for creating the data model.  
2-example.pdf Example graph.

## 3 Clustering Wikipedia Articles (R) (150 P)

In this task we form clusters of Wikipedia articles based on relative word frequencies. Data is provided in: /home/bigdata/8/enwiki-clean-10MiB.csv.

Using the k-means algorithm test separating the dataset into multiple clusters, inspect your results. Are the resulting clusters useful, do their members share characteristics?

Use hierarchical clustering on the same data, for an example, see: <https://www.r-bloggers.com/hierarchical-clustering-in-r-2/>

Discuss your programs execution time, compare the run time of hierarchical clustering to k-means.

### 3.1 Hints

For testing purposes working with a small subset of the data can be useful.

#### 3.1.1 Code Skeleton for R

```
1 library(tm)
2 d = read.csv("wikipedia-text-tiny-clean.csv", sep=";", quote="", header=F, stringsAsFactors=FALSE)
3 colnames(d) = c("title","content")
4
5 corpus = Corpus(VectorSource(d$content))
6
7 # We can set metadata, e.g. title:
8 i=0
9 corpus = tm_map(corpus, function(x) {
10   i = i + 1
11   meta(x, "title") = d$title[i]
12   x
13 })
14
15 # you can do:
16 # inspect(corpus)
17
18 # To access an element in the corpus use:
19 # corpus[[X]]
20
21 # we apply functions to clean the data
22 corpus = tm_map(corpus, removePunctuation)
23 corpus = tm_map(corpus, stripWhitespace)
24 corpus = tm_map(corpus, removeNumbers)
25 corpus = tm_map(corpus, content_transformer(tolower))
26 corpus = tm_map(corpus, removeWords, stopwords("english"))
27 # Create a sparse matrix for each document the word frequency
28 dtm = DocumentTermMatrix(corpus)
29 # To see data: e.g. call inspect(dtm[1:5,1:50] )
30
31 # Now apply clustering algorithms
32
33 # Inspect results....
```

### Submission:

3-wikipedia-clusters.pdf A lab notebook with your code, analysis and results.

## 4 HDFS REST API (Python) (120 P)

In this task, we explore the use of HDFS's REST API.

The specification of the API is available online: <http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/WebHDFS.html>.

For testing purpose, first, we use curl to perform the following operations:

- Listing of directory
- Creation of an file
- Fetching file content
- Deletion of a file

- Renaming a directory
- Deleting of a directory

Document your arguments to curl and the return values.

Now create a Python program which provides an FTP-alike interface to the user and performs the REST calls to WebHDFS, thus it allows to interactively browse, upload and download files from HDFS. The operations offered by the Python program should be analogous to the FTP <https://www.cs.colostate.edu/helpdocs/ftp.html>. Implement at least the following (hadoop fs) commands: ls, put, get, mkdir, rmdir and rm [-r]. Also implement the local operations ls and chdir.

## 4.1 Hints

If necessary, you may use <http://httpbin.org/> to debug your HTTP requests. The authentication on our HDFS is disabled, thus, you have to provide the query argument `user.name` with your name for all modifying operations.

### 4.1.1 Code skeleton

```

1  #!/usr/bin/env python3
2  import requests
3
4  def read_test():
5      # reading
6      # via url
7      url = 'http://10.0.0.61:50070/webhdfs/v1/user/gresens/data?op=OPEN'
8      r = requests.get(url)
9      print(r)
10     print(r.text)
11
12 def run_cmd(cmd):
13     """Translate command into REST API call and process the response."""
14     pass
15
16 def repl():
17     while True:
18         try:
19             cmd = input('>>> ')
20             except EOFError:
21                 print('exit')
22                 break
23             if cmd in ['exit', 'bye']:
24                 break
25             elif cmd == '':
26                 continue
27             else:
28                 run_cmd(cmd)
29
30 if __name__ == '__main__':
31     read_test()
32     repl()

```

### Submission:

- 4-curl.txt Your documented curl calls with responses.
- 4-hdfs-ftp.py Your Python FTP-alike HDFS REST client.

## 5 Optional bonus task: Finding places near you (60 P)

For this task, you will use the Google geocoding API to enable fuzzy finding of places from the wikipedia near you. The script should be invoked as follows: `./near.py Hamburg` and should in this case return a list of articles that are tagged as being located near Hamburg. Take a look at the documentation: <https://developers.google.com/maps/documentation/geocoding/intro>. Please consider that the API is subject to rate-limits, please don't access it repeatedly in an automated fashion.

### Submission:

- 5-near.py Your script to find places close to a given location.