

## S4 - WORDNET

### Assignment grading.

- Correctness of `SAP.java` (Mooshak score): 35%
- Correctness of `WordNet.java` (Mooshak score): 25%
- Correctness of the `Outcast.java` (Mooshak score): 20%
- Answers to questions (in `readme.txt`): 20%

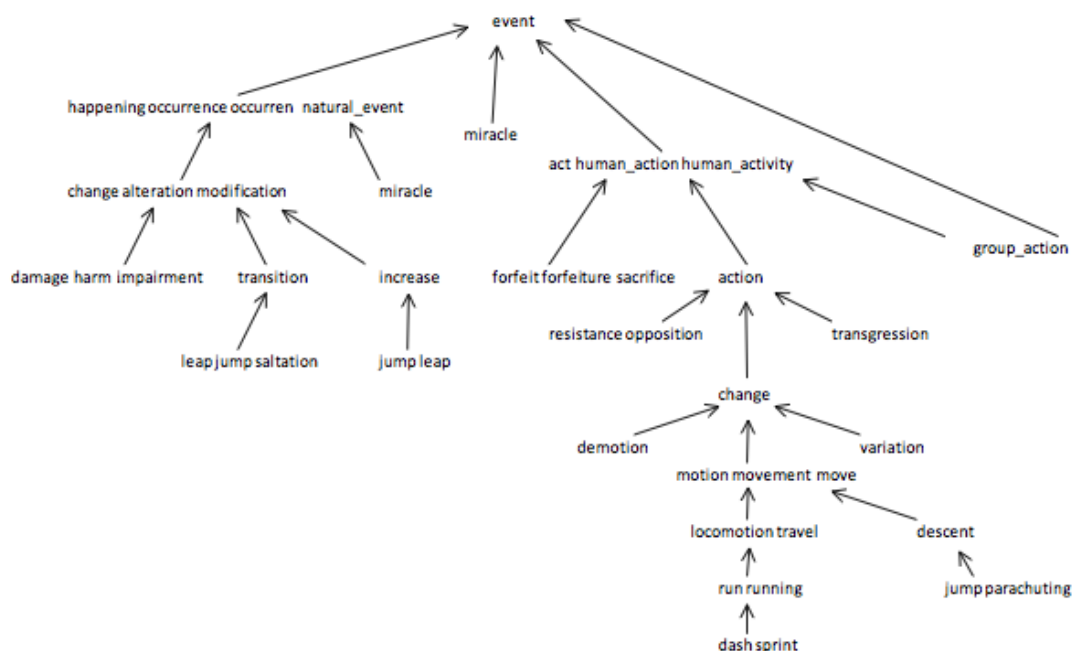
**Handin** Submit `WordNet.java`, `SAP.java` and `Outcast.java` on Mooshak. Submit `readme.txt` on MySchool.

**Note** You can submit each file to Mooshak at most 12 times. When working in pairs, this counts the submissions of both parties. Each additional submission reduces the grade by 0.05.

### INTRODUCTION

[WordNet](#) is a semantic lexicon for the English language that is used extensively by computational linguists and cognitive scientists; for example, it was a key component in IBM's [Watson](#). WordNet groups words into sets of synonyms called *synsets* and describes semantic relationships between them. One such relationship is the is-a relationship, which connects a *hyponym* („undirheiti“, more specific synset) to a *hypernym* („yfirheiti“, more general synset). For example, a *plant organ* is a hypernym of *carrot* and *plant organ* is a hypernym of *plant root*.

**The WordNet digraph** Your first task is to build the wordnet digraph: each vertex  $v$  is an integer that represents a synset, and each directed edge  $v \rightarrow w$  represents that  $w$  is a hypernym of  $v$ . The wordnet digraph is a *rooted DAG*: it is acyclic and has one vertex that is an ancestor of every other vertex. However, it is not necessarily a tree because a synset can have more than one hypernym. A small subgraph of the wordnet digraph is illustrated below.



**The WordNet input file formats** We now describe the two data files that you will use to create the wordnet digraph. The files are in *CSV format*: each line contains a sequence of fields, separated by commas.

- *List of noun synsets.* The file `synsets.txt` lists all the (noun) synsets in WordNet. The first field is the *synset id* (an integer), the second field is the synonym set (or *synset*), and the third field is its dictionary definition (or *gloss*). For example, the line

```
36,AND_circuit AND_gate,a circuit in a computer that fires only when all
of its inputs fire
```

means that the synset { `AND_circuit`, `AND_gate` } has an id number of 36 and its gloss is `a circuit in a computer that fires only when all of its inputs fire`. The individual nouns that comprise a synset are separated by spaces (and a synset element is not permitted to contain a space). The  $S$  synset ids are numbered 0 through  $S - 1$ ; the id numbers will appear consecutively in the synset file.

- *List of hypernyms* The file `hypernyms.txt` contains the hypernym relationships: The first field is a synset id; subsequent fields are the id numbers of the synset's hypernyms. For example, the following line

```
164,21012,56099
```

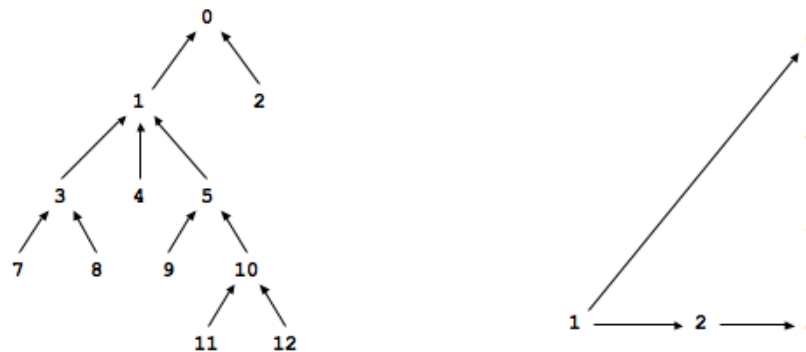
means that the the synset 164 (“`Actifed`”) has two hypernyms: 21012 (“`antihistamine`”) and 56099 (“`nasal_decongestant`”), representing that `Actifed` is both an antihistamine and a nasal decongestant. The synsets are obtained from the corresponding lines in the file `synsets.txt`.

```
164,Actifed,trade name for a drug containing an antihistamine and a
decongestant...
21012,antihistamine,a medicine used to treat allergies...
```

56099,nasal\_decongestant,a decongestant that provides temporary relief  
of nasal...

### SHORTEST ANCESTRAL PATH

An *ancestral path* between two vertices  $v$  and  $w$  in a digraph is a directed path from  $v$  to a common ancestor  $x$ , together with a directed path from  $w$  to the same ancestor  $x$ . A *shortest ancestral path* is an ancestral path of minimum total length. For example, in the digraph at left ([digraph1.txt](#)), the shortest ancestral path between 3 and 11 has length 4 (with common ancestor 1). In the digraph at right ([digraph2.txt](#)), one ancestral path between 1 and 5 has length 4 (with common ancestor 5), but the shortest ancestral path has length 2 (with common ancestor 0).



**SAP Data type** Implement an immutable data type SAP with the following API:

```

1 // constructor takes a digraph (not necessarily a DAG)
2 public SAP(Digraph G)
3
4 // length of shortest ancestral path between v and w; -1 if no such path
5 public int length(int v, int w)
6
7 // a common ancestor of v and w that participates in a shortest ancestral path; -1 if no such path
8 public int ancestor(int v, int w)
9
10 // length of shortest ancestral path between any vertex in v and any vertex in w; -1 if no such path
11 public int length(Iterable<Integer> v, Iterable<Integer> w)
12
13 // a common ancestor that participates in shortest ancestral path; -1 if no such path
14 public int ancestor(Iterable<Integer> v, Iterable<Integer> w)
15
16 // do unit testing of this class
17 public static void main(String[] args)

```

Your implementation must verify that the given digraph is actually a rooted DAG. All methods should throw a `java.lang.IndexOutOfBoundsException` if one (or more) of the input arguments is not between 0 and  $G.V()-1$ . You may assume that the iterable arguments contain at least one integer. All methods (and the constructor) should take time at most proportional to  $E + V$  in the worst case, where  $E$  and  $V$  are the number of edges and vertices in the digraph, respectively. Your data type should use space proportional to  $E + V$ .

**Test client** Your `main()` should take the name of a digraph input file as a command-line argument, construct the digraph, and use that digraph to test that all SAP methods work as expected. The following test client reads in vertex pairs from standard input, and prints out the length of the shortest ancestral path between the two vertices and a common ancestor that participates in that path: On the left below you can see the contents of `digraph1.txt` and on the right are a few samples of expected values if your program reads in that input file:

13		<code>v = 3, w = 11</code>
11		<code>length = 4, ancestor = 1</code>
7	3	
8	3	<code>v = 9, w = 12</code>
3	1	<code>length = 3, ancestor = 5</code>
4	1	
5	1	<code>v = 1, w = 6</code>
9	5	<code>length = -1, ancestor = -1</code>
10	5	
11	10	<code>v = {3, 9, 7, 1}, w = {11, 2, 6}</code>
12	10	<code>length = 2, ancestor = 0</code>
1	0	
2	0	

## WORDNET DATA TYPE

Implement an immutable data type `WordNet` with the following API:

```

1 // constructor takes the name of the two input files
2 public WordNet(String synsets, String hypernyms)
3
4 // returns all WordNet nouns
5 public Iterable<String> nouns()
6
7 // is the word a WordNet noun?
8 public boolean isNoun(String word)
9
10 // distance between nounA and nounB (defined below)
11 public int distance(String nounA, String nounB)
12
13 // a synset (second field of synsets.txt) that is the common
14 ancestor of nounA and nounB
15 // in a shortest ancestral path (defined below)
16 public String sap(String nounA, String nounB)
17
18 // do unit testing of this class
19 public static void main(String[] args)

```

The constructor should throw a `java.lang.IllegalArgumentException` if the input does not correspond to a rooted DAG. The `distance()` and `sap()` methods should throw a `java.lang.IllegalArgumentException` unless both of the noun arguments are `WordNet` nouns.

Your data type should use space linear in the input size (size of synsets and hypernyms files). The constructor should take time linearithmic (or better) in the input size. The method `isNoun()` should run in time

logarithmic (or better) in the number of nouns. The methods `distance()` and `sap()` should run in time linear in the size of the WordNet digraph.

## MEASURING SEMANTIC RELATEDNESS OF TWO NOUNS

Semantic relatedness refers to the degree to which two concepts are related. Measuring semantic relatedness is a challenging problem. For example, most of us agree that *George Bush* and *John Kennedy* (two U.S. presidents) are more related than are *George Bush* and *chimpanzee* (two primates). However, not most of us agree that *George Bush* and *Eric Arthur Blair* are related concepts. But if one is aware that *George Bush* and *Eric Arthur Blair* (aka George Orwell) are both communicators, then it becomes clear that the two concepts might be related.

**Outcast detection** Given a list of wordnet nouns  $A_1, A_2, \dots, A_n$ , which noun is the least related to the others? To identify an outcast, compute the sum of the distances between each noun and every other one:

- $distance(A, B)$  = distance is the minimum length of any ancestral path between any synset  $v$  of  $A$  and any synset  $w$  of  $B$ .

Implement an immutable data type `Outcast` with the following API:

```

1 // constructor takes a WordNet object
2 public Outcast(WordNet wordnet)
3
4 // given an array of WordNet nouns, return an outcast
5 public String outcast(String[] nouns)
```

Assume that argument array to the `outcast()` method contains only valid wordnet nouns (and that it contains at least two such nouns).

The following test client takes from the command line the name of a synset file, the name of a hypernym file, followed by the names of outcast files, and prints out an outcast in each file:

```

1 public static void main(String[] args) {
2     WordNet wordnet = new WordNet(args[0], args[1]);
3     Outcast outcast = new Outcast(wordnet);
4     for (int t = 2; t < args.length; t++) {
5         String[] nouns = In.readStrings(args[t]);
6         StdOut.println(args[t] + ": " + outcast.outcast(nouns));
7     }
8 }
```

SCHOOL OF COMPUTER SCIENCE, REYKJAVÍK UNIVERSITY, MENNTAVEGI 1, 101 REYKJAVÍK

E-mail address: mmh@ru.is