

“user/comp9313/output”. The input and output paths are obtained from the arguments.

Question 2. Download the sample input file “Votes.csv” from: <https://webcms3.cse.unsw.edu.au/COMP9313/17s2/resources/12622>, and put it in HDFS folder “/user/comp9313/input”. In this file, the fields are separated by ‘,’ and the lines are separated by ‘\n’. The data format of “Votes.csv” is as below:

```
- Id
- PostId
- VoteTypeId
  - ` 1`: AcceptedByOriginator
  - ` 2`: UpMod
  - ` 3`: DownMod
  - ` 4`: Offensive
  - ` 5`: Favorite - if VoteTypeId = 5 UserId will be populated
  - ` 6`: Close
  - ` 7`: Reopen
  - ` 8`: BountyStart
  - ` 9`: BountyClose
  - `10`: Deletion
  - `11`: Undeletion
  - `12`: Spam
  - `13`: InformModerator
  - `14`:
  - `15`:
  - `16`:
- UserId (only for VoteTypeId 5)
- CreationDate
```

(i). Find the top-5 VoteTypeIds that have the most distinct posts. You need to output the VoteTypeId and the number of posts. The results are ranked in descending order according to the number of posts, and each line is in format of: VoteTypeId\tNumber of posts.

(ii). Find all posts that are favoured by more than 10 users. You need to output both PostId and the list of UserIds, and each line is in format of:

PostId#UserId1,UserId2,UserId3,...,UserIdn

The lines are sorted according to the NUMERIC values of the PostIds in ascending order. Within each line, the UserIds are sorted according to their NUMERIC values in ascending order.

(Hint: the mkString function is useful to format your output)

You can download the code template at:

<https://webcms3.cse.unsw.edu.au/COMP9313/17s2/resources/12624>.

You can see the result at:

<https://webcms3.cse.unsw.edu.au/COMP9313/17s2/resources/12623>.