

```
1 from mrjob.job import MRJob
2 import re
3
4 # Regular expression to split words by non-alphabetic
  characters
5 WORD_RE = re.compile(r"\b\w+\b")
6
7
8 class InvertedIndex(MRJob):
9
10     def mapper(self, _, line):
11         # Check if the line has the expected format
12         if ":" in line:
13             document_number, text = line.split(':', 1
14         )
15             words = re.findall(WORD_RE, text.lower())
16             for word in words:
17                 yield (word, document_number.strip())
18             else:
19                 # Handle lines with unexpected format (e.
20                 g., empty lines)
21                 pass
22
23     def reducer(self, word, documents):
24         # Create a set to store unique document IDs
25         unique_documents = set()
26
27         # Add each document ID to the set
28         for document in documents:
29             unique_documents.add(document)
30
31         # Emit the word and the list of documents
32         where it appears
33         yield (word, sorted(list(unique_documents)))
34
35 if __name__ == '__main__':
36     InvertedIndex.run()
```