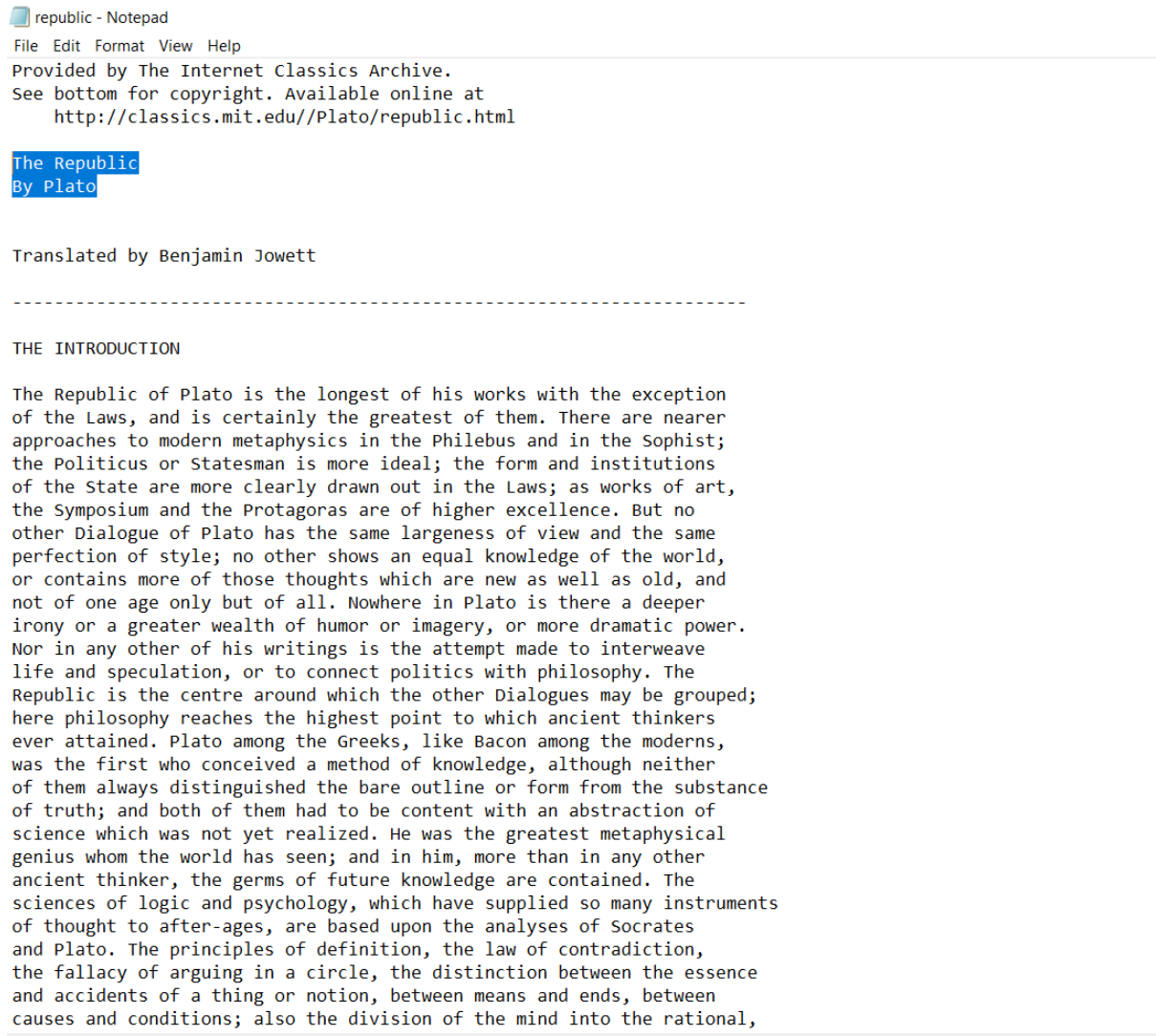Yogeshwaran S

Case Study 8

Module 8: Data Flow and Data Processing


Step 1 –

I create a txt file of a book called "The Republic"
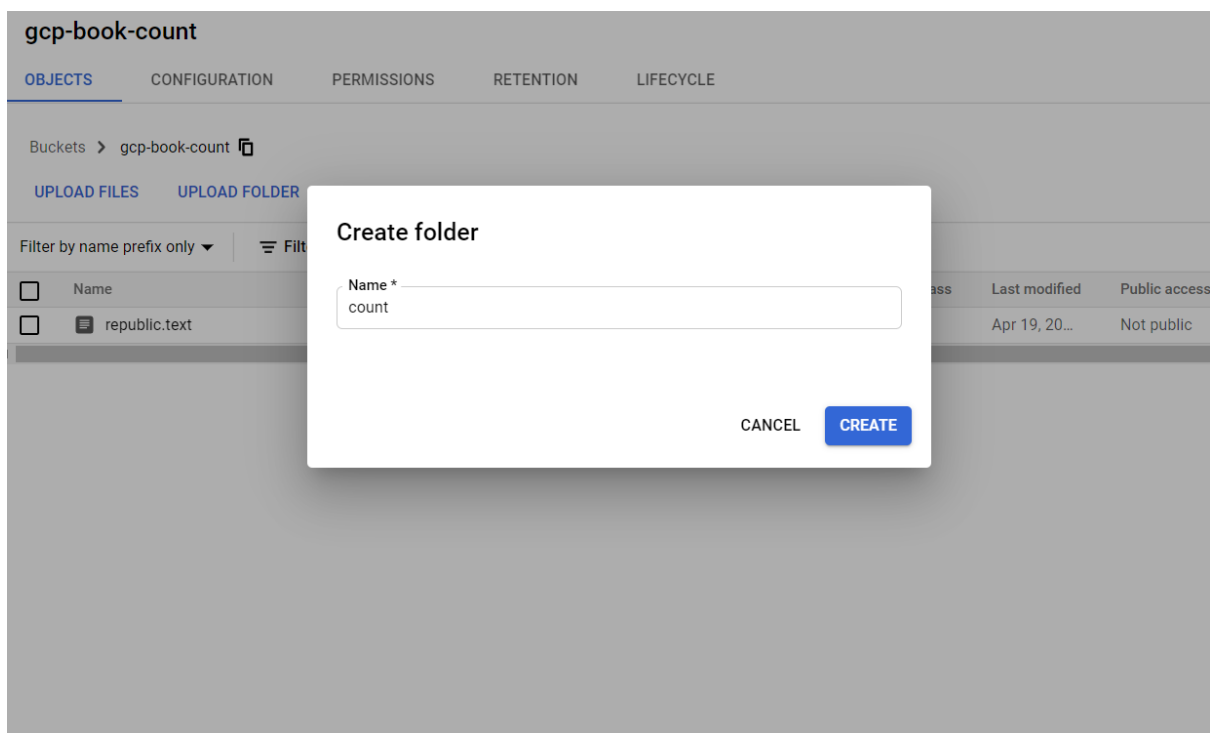
republic - Notepad

File  Edit  Format  View  Help

Provided by The Internet Classics Archive.
See bottom for copyright. Available online at
    http://classics.mit.edu//Plato/republic.html

The Republic
By Plato


Translated by Benjamin Jowett

-------------------------------------------------------------------

THE INTRODUCTION

The Republic of Plato is the longest of his works with the exception
of the Laws, and is certainly the greatest of them. There are nearer
approaches to modern metaphysics in the Philebus and in the Sophist;
the Politicus or Statesman is more ideal; the form and institutions
of the State are more clearly drawn out in the Laws; as works of art,
the Symposium and the Protagoras are of higher excellence. But no
other Dialogue of Plato has the same largeness of view and the same
perfection of style; no other shows an equal knowledge of the world,
or contains more of those thoughts which are new as well as old, and
not of one age only but of all. Nowhere in Plato is there a deeper
irony or a greater wealth of humor or imagery, or more dramatic power.
Nor in any other of his writings is the attempt made to interweave
life and speculation, or to connect politics with philosophy. The
Republic is the centre around which the other Dialogues may be grouped;
here philosophy reaches the highest point to which ancient thinkers
ever attained. Plato among the Greeks, like Bacon among the moderns,
was the first who conceived a method of knowledge, although neither
of them always distinguished the bare outline or form from the substance
of truth; and both of them had to be content with an abstraction of
science which was not yet realized. He was the greatest metaphysical
genius whom the world has seen; and in him, more than in any other
ancient thinker, the germs of future knowledge are contained. The
sciences of logic and psychology, which have supplied so many instruments
of thought to after-ages, are based upon the analyses of Socrates
and Plato. The principles of definition, the law of contradiction,
the fallacy of arguing in a circle, the distinction between the essence
and accidents of a thing or notion, between means and ends, between
causes and conditions; also the division of the mind into the rational,


Step 2 create a bucket called "book" on gcp and move this txt file

## Step 3 create a output folder and temp folder inside the bucket



## Step 4 create a dataflow

I select word count template from dataflow

Step 5

I fill necessary things

## Required parameters

**Input file(s) in Cloud Storage \***

gs://gcp-book-count/republic.text

The input file pattern Dataflow reads from. Use the example file (gs://dataflow-samples/shakespeare/kinglear.txt) or enter the path to your own using the same format: gs://your-bucket/your-file.txt

**Output Cloud Storage file prefix \***

gs://gcp-book-count/count

Path and filename prefix for writing output files. Ex: gs://your-bucket/counts

**Temporary location \***

gs://gcp-book-count/temp

Path and filename prefix for writing temporary files. Ex: gs://your-bucket/temp

## Encryption

🔘 Google-managed encryption key
No configuration required

⭘ Customer-managed encryption key (CMEK)
Manage via Google Cloud Key Management Service

⌄ SHOW OPTIONAL PARAMETERS

**RUN JOB**

Step 6

Job ready

## Step 7

Now I can see the outputs in my bucket



## Step 8

Output is here and I will attach my output file also from this report