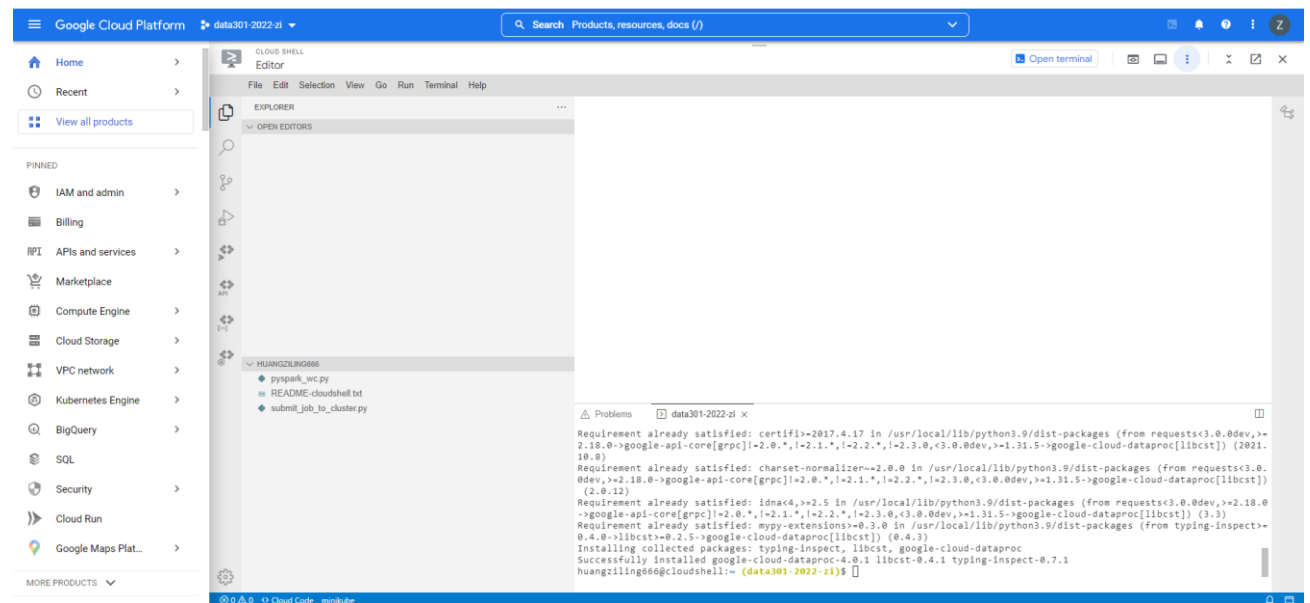Name: Ziling Huang

Username: zhu51
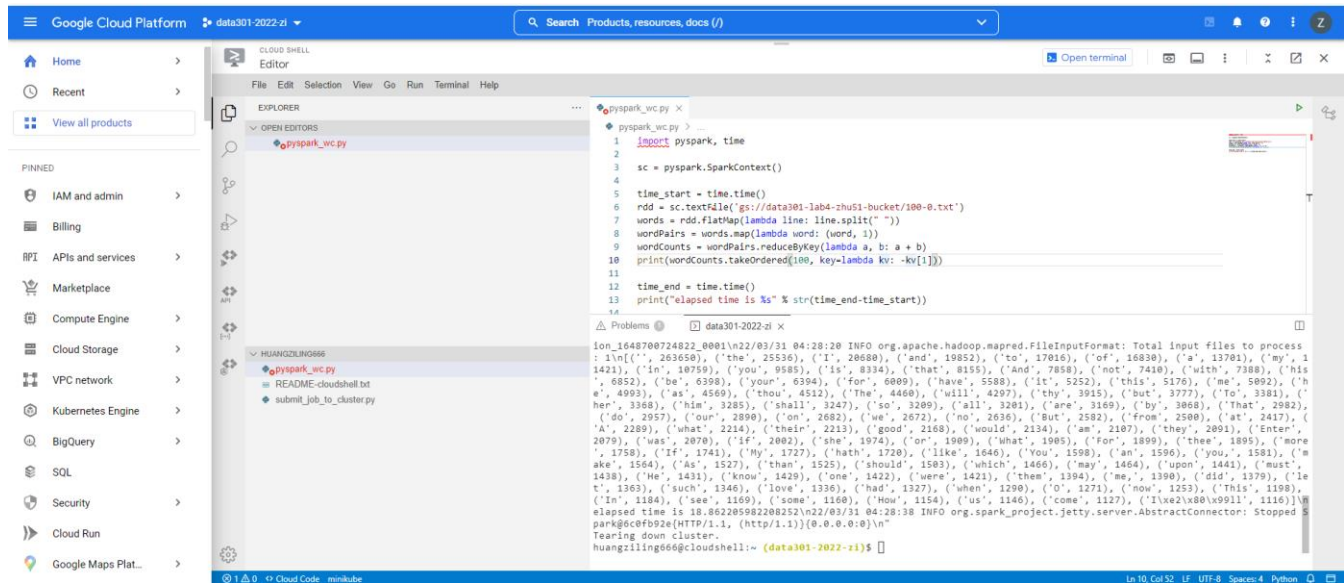
## Part2



## Part3

REGION=australia-southeast1
ZONE=australia-southeast1-a
PROJECT=**data301-2022-zi**
CLUSTER=data301-lab4-**zhu51**-cluster
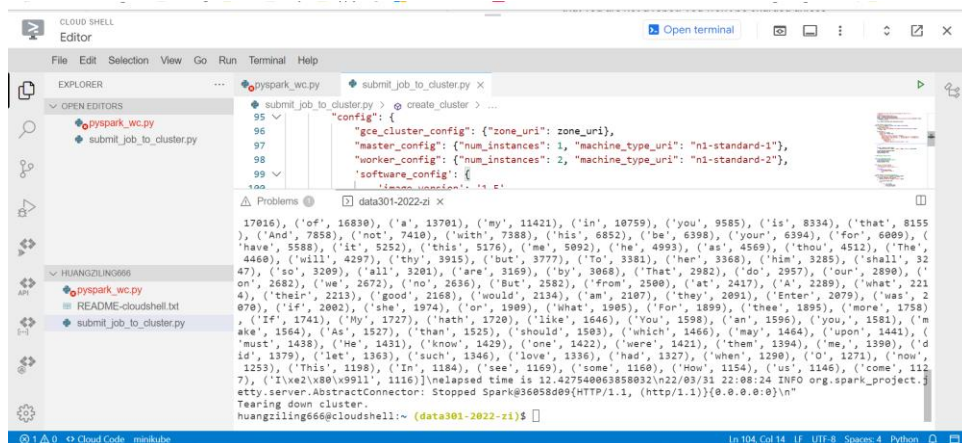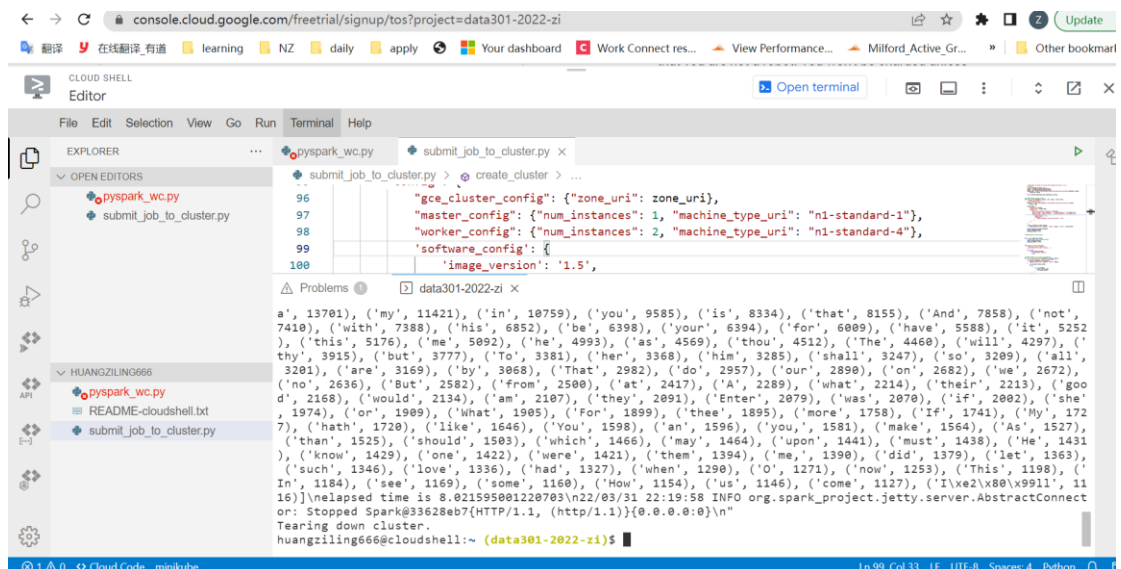BUCKET=data301-lab4-**zhu51**-bucket
When 2 cores:

# Part4

4b)

4 cores:



8 cores:

16 cores:



4c)

One core:



4d)

| Tseq | | |
|---|---|---|
| 21.02418566 | | |
| | | |
| | Tpar | speedup |
| p = 2 | 18.86220598 | 1.11462 |
| p = 4 | 12.42754006 | 1.691742 |
| p = 8 | 8.021595001 | 2.620948 |
| p = 16 | 7.503222942 | 2.802021 |

4e)

As we add more processors, the processing time is decreasing as well as the speedup time is increasing, which means the more processors we use, the less time we spend. However, if we the add processors more than 16, the speedup time and the processing time will be the same since the ability has already reached a limit so more processors won't increase the efficiency anymore.

4f)