# Assignment 4

## Task1

```
The  frequency  position  of  the  words  'applicant'  is:    448
The  frequency  position  of  the  words  'and'  is:    2
The  frequency  position  of  the  words  'attack'  is:    512
The  frequency  position  of  the  words  'protein'  is:    3167
The  frequency  position  of  the  words  'car'  is:    648
```

my-project-zipengxu > as4

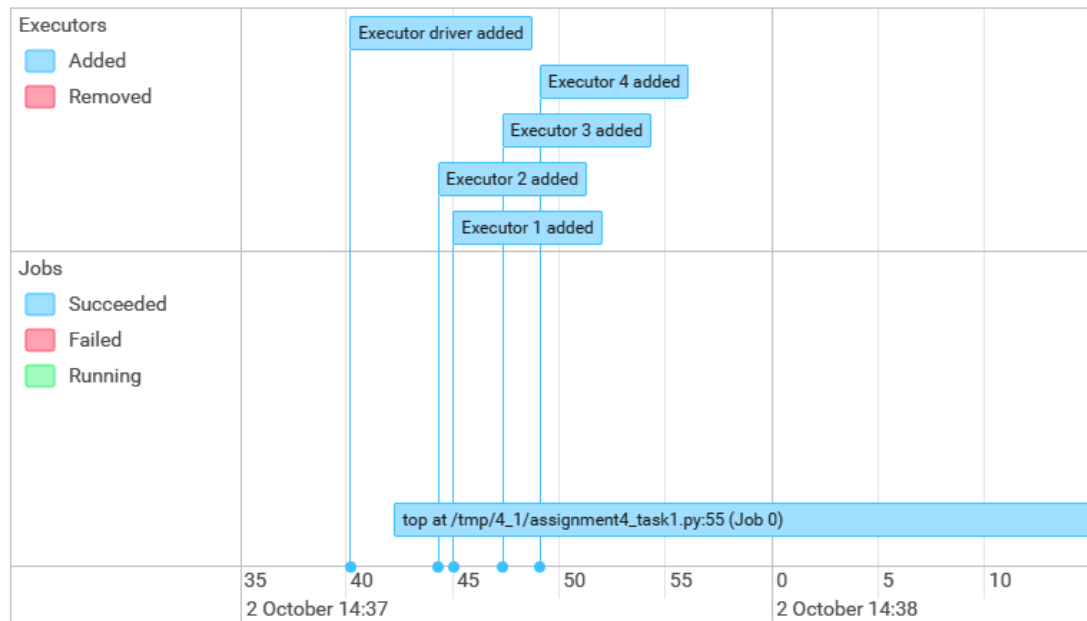Spark 3.1.2 | Jobs | Stages | Storage | Environment | Executors

### Spark Jobs (?)

**User:** root
**Total Uptime:** 1.9 min
**Scheduling Mode:** FAIR
**Completed Jobs:** 7

▼ Event Timeline
☐ Enable zooming

Executors
- ☐ Added
- ☐ Removed

Executor driver added
Executor 4 added
Executor 3 added
Executor 2 added
Executor 1 added

Jobs
- ☐ Succeeded
- ☐ Failed
- ☐ Running

top at /tmp/4_1/assignment4_task1.py:55 (Job 0)

35    40    45    50    55    0    5    10
2 October 14:37                    2 October 14:38

# Task2



Top 5 words:
['that', 'not', 'tribunal', 'court', 'applicant']

# Task3





F1-score of my classifier:
0.0394743730694728

I sample 3 articles that my model though were Australian court case (id are :['7524908', '7522559', '7522492']). The reasons that my model was fooled are as followed:

1. I pick these words which are highly correlative with courts articles are basically included in legal or court cases article, which means it is unlikely to take 'Australian' court cases article precisely from the huge Wikipedia articles dataset. If I weight some words higher,

like 'Australian', in my model, I may be more likely to pick Australian court cases precesly.
2. The learning rate I select is based on the training model using small unbalance dataset. I might take a relatively large learning rate in my model, which means I fail to get the local minimum variable in gradient descent. So my model could be improved by adapting parameters.