

PP4 Report

Logan Bontrager

November 10, 2020

Task 1)

The top five words for each topic are shown below where each line corresponds to a topic and the words are presented in decreasing order.

```
1  earth,question,large,temperature,things
2  don,driving,shift,cars,mph
3  space,science,internet,world,information
4  station,shuttle,option,launch,redesign
5  mission,hst,pat,solar,shuttle
6  car,clutch,shifter,drive,gear
7  edu,gif,insurance,uci,ics
8  car,ford,probe,seat,nice
9  cars,heard,diesels,people,matter
10 edu,writes,article,apr,find
11 space,nasa,gov,such,long
12 back,used,small,sho,use
13 oil,time,service,come,change
14 sky,bill,moon,light,rights
15 car,book,mustang,rear,price
16 system,oort,edu,david,ray
17 cost,second,years,cars,comes
18 engine,turbo,power,steering,temp
19 henry,edu,toronto,spencer,writes
20 make,don,point,want,good
```

Here, we can see that most topics have well defined themes. It appears that the two underlying topics present in the data are cars and space. One thing that sticks out to me about the model is that some words carry more weight

in terms of human theme definition than others. For example, topic one contains the word things which doesn't have much weight in terms of personal interpretability although I suppose it could be useful in other ways.

Task 2)

The graphs below show the test error vs. training size for the bayesian logistic regression model trained on both the topic representation data and the bag of words. we can see that both datasets performance similarly. The main distinction is that the topic representation data appears to perform better on smaller training sizes than the bag of words. This could be advantageous under data scarcity. The error bars from each graph look rather consistent with each other given the y scale.

