

# CS5740: Assignment 1

https:

[//github.com/cornell-cs5740-21sp/a1--group-5](https://github.com/cornell-cs5740-21sp/a1--group-5)

Kae-Jer (Mike) Cho  
kc2225

Vianne Gao  
vrg36

Yuxin Zhang  
yz2726

## 1 Introduction (5pt)

We implement the Perceptron model and the multi-layer Perceptron (MLP) model and apply them to two tasks: Proper Name Classification where we take proper names and predict whether they indicate people, movie, drugs, place and or company; and Newsgroup Classification where we classify documents into 20 topics.

We experiment with different input features for each dataset and model separately, perform ablation tests on features to identify the most informative inputs, and experiment with different model complexity, activation functions and optimizers for the MLP model. Overall, we observed that MLP outperforms Perceptron in both tasks, and both models are generalizable to unseen dataset.

## 2 Features (20pt)

There are three feature components in both propername data feature and newsgroup data feature when doing feature extraction (N-gram, Task-specific features, and combined features).

For the task-specific feature for newsgroup dataset, we engineered 4 features based on the following observation:

- Each group is distinguished by whether headers such as NNTP-Posting-Host: and Distribution: appear more or less often. Feature **host\_count**: count the num. of "NNTP-Posting-Host:". Feature **distribution\_count**: count the num. of "Distribution:".
- Another significant feature involves whether the sender is affiliated with a university, as indicated either by their headers or their signature. Feature **from\_university**: binary variable for whether email address of sender contains "edu".
- The word "article" is a significant feature based on how often people quote previous posts like: "In article [article ID], [name] <[e-mail address]> wrote:". Feature **reference\_count**: count the num. of "In article" in text.

For the task-specific feature for propername dataset, we engineered 4 features based on the following observation:

- Text containing numbers can hardly be names. Feature **contain\_numerics**: binary variable whether text contains numbers.

- Text containing some punctuation (-: "%) can hardly be names. Feature **contain\_special\_punc**: binary variable whether text contains special punctuation (&: "%).
- Text containing "Inc." can be strong indicator for the class company. Feature **contain\_inc**: binary variable whether text contains "Inc." in the end.
- Text containing only one token can hardly be a proper name. Feature **small\_token\_length**: binary variable whether text contains only token length 1.

	Propername dataset	Newsgroup dataset
N-gram	Character-level N-gram feature engineering, controlled by N and MinFreq(see text)	Bag of words model Feature engineering controlled by N and MinFreq(see text)
Task-specific	See feature definitions above	See feature definitions above
Combined	Horizontally stack N-gram features and task-specific features	Horizontally stack N-gram features and task-specific features

Table 1: Feature components

The word and character frequency dictionary is built on the text from the entire dataset to maximize the inclusion of features and avoid unknown words or characters. In our N-gram models, we allow input of N to be a tuple specifying what range of word or character sequence should the feature engineering model consider. MinFreq (minimal Frequency) is introduced for the feature engineering process to select the most prominently appeared word or character, thus reducing dimensionalities.

Table 2 shows an example of how the generated three feature components look like on the propername dataset with  $N = (1,1)$  and  $\text{MinFreq} = 15$ .

## 3 Experimental Setup

**Data (5pt)** We divide the Proper Name dataset into training, development and test sets. The



Model (Newsgroup)	Accuracy
<b>Development Results</b>	
Perceptron (Comb., $N=(1,3)$ )	<b>0.8714</b>
Perceptron w/o Task-Specific Feat.	0.8608
Perceptron w/o 3-Gram	0.8612
Perceptron w/o 2,3-Gram	0.8608
Perceptron w/o 1-Gram	0.6195
MLP (Comb., $N=(1,2)$ , sigmoid,Adam)	<b>0.9152</b>
MLP (sigmoid,AdaGrad)	0.9089
MLP (tanh,Adam)	0.9046
MLP w/o 2-Gram	0.8984
MLP w/o 1-Gram	0.7548
<b>Test Results</b>	
Perceptron	0.7521
MLP	0.8111

Table 6: Performance on 20Newsgroup. We evaluate the performance of the model when certain features are ablated and when different activation functions and optimizers are used. We mark in bold the best results. \*Comb. indicates combined features.

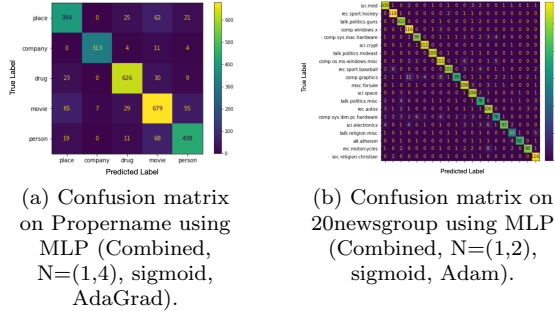


Figure 1: Confusion Matrices

We performed feature and parameter search separately for each model. Our final Perceptron model is selected using the development set. We experiment with different combinations of  $N$ 's for the  $N$ -gram feature, with and without the task-specific features. We observed a minor drop in performance when excluding the task-specific features and 4-Grams, but see bigger drops in accuracy when 1-Gram and 3-Gram are ablated, highlighting their importance in the classification task. We can also find the similar drop from the MLP model when 1-Gram are ablated.

From Figure 1a, we see that our MLP model is prone to error when a word may belong to either class under different context. For example, the name of a movie can also be the name of a 'person' (e.g. Queen Victoria) or of a 'place' (e.g. Peters). Such ambiguity makes it difficult for the model to distinguish these related classes.

## 4.2 Newsgroup Classification (12pt)

On the test set, we achieved an accuracy of 0.7521 using Perceptron, and the accuracy is improved

Batch Size	Average	S.D.
<b>Proper Name</b>		
1	333.17	2.81
300	17.06	0.06
500	16.56	0.44
1000	16.29	0.56
<b>Newsgroup</b>		
1	20.93	0.14
300	0.75	0.01
500	0.71	0.01
1000	1.35	0.53

Table 7: Batching Benchmarking

to 0.8111 using MLP. The input features to both models include  $N$ -gram for  $N \in 1, 2$  and the task-specific features. We also included 3-Gram for Perceptron.

We experiment with different combinations of  $N$ 's for the  $N$ -gram feature, with and without the task-specific features. We observed a minor drop in performance when excluding the task-specific features, 2 and 3-Grams, but see bigger drops in accuracy when 1-Gram is ablated, high-lighting the importance of 1-Gram task with both Perceptron and MLP model.

From Figure 1b, we see the MLP model is prone to error when two classes of documents are related and share common key words. For example, the topic 'comp.graphics' is highly related to 'comp.windows.x', hence a lot of vocabularies and phrases are shared.

## 4.3 Batching Benchmarking (4pt)

The result of the batching benchmarking of MLP model on GPU is shown in Table 7. The speed is faster with batching than without. The reason is that using larger batch sizes would allow us to parallelize computations to a greater degree, hence speeding up the computation time. The training time of the Newsgroup is faster than the Proper Name in our result. The reason is that the epoch number of each model is different (11 for Proper Name and 200 for Newsgroup), so even the feature dimensionality of the Newsgroup is larger than the Proper Name, the training speed for Newsgroup's model is faster.

## 5 Conclusion (4pt)

Our Perceptron and MLP models are fairly generalizable and achieve reasonable accuracy on the 2 classification tasks. We found that 1-Gram is important for the performance in both datasets, and MLP consistently outperforms Perceptron. Qualitative analysis of our models show that most errors are made when distinguishing between ambiguous names and between related topics with the same theme.