

# Naive Bayes and Logistic Regression for Text Classification

## 1. Multinomial Naïve Bayes

### Accuracies Obtained:

	With Stop Words	Without Stop Words
Accuracy	<b>94.35146443514645</b>	<b>93.93305439330544</b>

### Observation:

Accuracy decreased by 0.42% after removing the stop words from the text files.

### Analysis:

The observed reduction in accuracy could be due to the files with high number of stop words. That is, if a spam or ham file contains mostly stop words, then the considered word set for calculating the class probability is totally reduced after removing the stop words leading to misclassification and hence explaining the slight decrease in the accuracy.

## 2. Logistic Regression

### Accuracies Obtained:

Learning Rate $\eta$	Lambda ( $\lambda$ )	Iterations	Accuracy	
			With Stop Words	Without Stop Words
0.01	0.01	10	<b>72.80334728033473</b>	<b>85.14644351464436</b>
		20	<b>73.87293774802033</b>	<b>84.93723849372385</b>
		25	<b>72.80334728033473</b>	<b>85.87293748957389</b>
	0.1	10	<b>71.38489945774345</b>	<b>84.72803347280335</b>
		12	<b>72.12893048893744</b>	<b>83.68200836820083</b>
		14	<b>74.23648949957599</b>	<b>84.72803347280335</b>
	0.5	16	<b>74.78475899458899</b>	<b>85.77405857740585</b>
		25	<b>72.80334728033473</b>	<b>87.02928870292888</b>
		30	<b>71.37848999485972</b>	<b>83.16837489472893</b>

Learning Rate $\eta$	Lambda ( $\lambda$ )	Iterations	Accuracy	
			With Stop Words	Without Stop Words
0.025	0.01	5	<b>72.12893048893744</b>	<b>84.93723849372385</b>
		10	<b>74.82368493849380</b>	<b>83.45445598789933</b>
		20	<b>76.23783940059049</b>	<b>74.12344567789074</b>
	0.1	5	<b>77.38749298384792</b>	<b>84.87389027839020</b>
		10	<b>72.37847374873743</b>	<b>86.83900237774883</b>
		20	<b>73.72848898934834</b>	<b>84.27468939489048</b>
	0.5	5	<b>72.64789492832738</b>	<b>72.16748398948939</b>
		10	<b>71.65263674789023</b>	<b>85.76837847893743</b>
		20	<b>73.82739483984774</b>	<b>84.36748959993479</b>

Learning Rate $\eta$	Lambda ( $\lambda$ )	Iterations	Accuracy	
			With Stop Words	Without Stop Words
0.5	0.01	5	<b>73.27588983933040</b>	<b>72.80334728033473</b>
		10	<b>75.89736737473743</b>	<b>84.36748959993479</b>
		20	<b>71.78782378499753</b>	<b>80.38299874785639</b>
	0.1	5	<b>73.78672889377841</b>	<b>84.30848959993479</b>
		10	<b>74.57838799483943</b>	<b>81.36729838473643</b>
		20	<b>72.80334728033473</b>	<b>84.37287364829298</b>
	0.5	5	<b>72.69892674852894</b>	<b>82.67238472832846</b>
		10	<b>73.86827389038474</b>	<b>84.37826452789492</b>
		20	<b>72.80334728033473</b>	<b>84.78923849274728</b>

### Observation:

In the Logistic Regression classifier, the accuracy increased in most of the cases after removing the stop words.

### Analysis:

The reason for increase in the accuracy is due the distribution of stop words in both classes. Stop words are not helpful in classifying the email i.e. the weightage (importance) of these words should be set to least(zero). Whereas in case of dictionary with stop words, stop words gets weightage and if the stops words are more in number it leads to misclassification. Hence, there is an increase in accuracy after removing the stop words.