

# Machine Learning: Final Project

Noah Buchanan

Zachary Mackay

Jacob Burton

Computer Science Department  
University of Arkansas - Fort Smith

May 9, 2021

## 1 Introduction

This was an experimentation process by us to determine the best hyper-parameters for different machine learning algorithms and the interaction of those hyper-parameters potentially if any. The Machine Learning algorithms chosen were Multi-layer Perceptron(MLP), Support Vector Machine(SVM), and AdaBoost.

## 2 Datasets

### Spam Dataset

Our first dataset was a dataset of spam or otherwise denoting so with a binary classification. The features include 48 word frequency features out of the total amount of words. 6 character frequency features out of the total amount of characters. Lastly capital\_run\_length\_average, capital\_run\_length\_longest, and capital\_run\_length\_total. Capital run in this case being a sting of capital letters bounded by a non alphanumeric character, and those denote the average longest and total lengths. Below are some examples of what the data looks like, the data is split so the first 1800 records are all spam and past that all not spam records but we randomly split the data so this is not important and the data has a ratio of about 40% spam to 60% not spam. Lastly all data was normalized to be on the same scale with z-scores.

[illegible]

## League of Legends Dataset

Our second dataset was various game statistics from the 10 minute mark into a match for a game called League of Legends. The features include things such as team total kills, team total deaths, objectives taken by either team, the gold(in-game money) differential between the teams. Our y column prediction was a binary class denoting 1 for a blue win and 0 for a red win. This data was once again normalized by z-score. I will provide examples of what the data looks like again below. The game ID was unimportant and dropped, it is used by riot game's api to uniquely identify matches and not important.

gameId	blueWins	blueWard	blueWard	blueFirst	blueKills	blueDeaths	blueAssis	blueElite	blueDrag	blueHeral	blueTower	blueTotal	blueAvgL	blueTotal	blueTotal	blueGoldD	blueExpe	blueCpe	blueGold	redWards	redWards	redFirstBl
4.52E+09	0	28	2	1	9	6	11	0	0	0	17210	6.6	17055	1955	36	-643	-8	19.5	1721	15	6	0
4.52E+09	0	12	1	0	5	5	5	0	0	0	14712	6.6	16265	174	43	-2908	-1173	17.4	1471.2	12	1	1
4.52E+09	0	15	1	0	7	11	4	1	1	0	16113	6.4	16221	186	46	-1172	-1033	18.6	1611.3	15	3	1
4.52E+09	0	43	1	0	4	5	5	1	0	1	13157	7	17954	201	55	-1321	-7	20.1	1315.7	15	2	1
4.48E+09	0	75	4	0	6	6	6	0	0	0	16600	7	18543	210	57	-1004	230	21	1640	17	2	1
4.48E+09	1	18	0	0	5	3	6	1	1	0	15899	7	18161	225	42	698	101	22.5	1589.9	36	5	1
4.49E+09	0	18	3	1	7	6	7	1	1	0	16874	6.8	16967	225	53	2411	1563	22.5	1687.4	57	1	0
4.5E+09	0	16	2	0	5	13	3	0	0	0	13305	6.4	16138	209	48	-2615	-800	20.9	1330.5	15	0	1
4.48E+09	0	16	3	0	7	7	8	0	0	0	16601	7.2	18527	189	61	-1979	-771	18.9	1660.1	15	2	1
4.51E+09	1	13	1	1	4	5	5	1	1	0	15057	6.8	16805	220	39	-1548	-1574	22	1505.7	16	2	0
4.45E+09	0	20	3	1	4	4	6	0	0	0	15474	6.6	16611	231	28	331	-1585	23.1	1547.4	15	2	0
4.45E+09	0	33	2	1	11	11	7	1	0	1	16695	7	18507	157	40	-1505	-635	15.7	1669.5	17	1	0
4.52E+09	1	18	1	1	7	1	11	1	1	0	17865	7.4	19102	236	53	3274	1659	23.8	1786.5	12	1	0
4.52E+09	0	14	4	0	9	1	1	1	0	1	14979	6.6	17213	210	52	-3414	-1141	21	1497.9	20	3	1
4.52E+09	1	15	3	1	4	4	4	0	0	0	15722	6.8	17896	224	51	-470	-187	22.4	1572.2	102	1	0
4.48E+09	0	17	1	0	3	7	3	0	0	0	15015	6.8	16974	209	53	-1996	-1804	20.9	1501.5	18	3	1
4.52E+09	1	14	1	1	10	2	8	0	0	0	19733	7.6	20862	263	56	5228	3978	26.3	1973.3	13	2	0
4.5E+09	0	43	3	0	3	7	3	1	0	1	14852	6.8	16888	203	54	-1975	-1345	20.3	1485.2	17	14	1
4.49E+09	1	21	4	1	5	4	11	0	0	0	16282	6.8	17378	213	49	882	512	21.3	1628.2	19	3	0
4.48E+09	0	11	3	0	5	9	5	0	0	0	14994	7	17924	188	48	-3155	-2773	18.8	1499.4	15	1	1
4.47E+09	1	14	3	1	11	6	15	1	1	0	18606	7.2	19254	214	45	2536	1625	21.4	1860.6	12	2	0
4.45E+09	0	13	1	0	4	13	5	0	0	0	15878	6.8	16945	203	43	-2835	-2821	20.3	1587.8	15	3	1
4.48E+09	0	17	2	0	4	6	3	0	0	0	15773	7	17562	238	56	-1254	-567	23.8	1577.3	16	3	1
4.52E+09	0	78	4	0	4	3	4	2	1	1	13906	6.8	19219	238	67	406	1283	23.8	1390.6	15	3	1
4.52E+09	1	13	3	0	6	9	10	0	0	0	16155	6.8	17223	208	45	-1020	493	20.8	1615.5	18	3	1
4.52E+09	0	17	1	0	4	8	4	0	0	0	15415	6.2	16047	204	44	-2472	-1067	20.4	1541.5	19	2	1
4.51E+09	1	18	1	0	8	6	9	1	1	0	17178	7	17854	209	44	1260	-55	20.9	1717.8	10	3	1
4.5E+09	0	19	3	1	7	4	9	0	0	0	17226	7	18982	228	61	1857	2015	22.8	1722.6	14	4	0
4.52E+09	0	105	5	1	8	6	13	1	1	0	17474	6.8	17486	217	53	1120	-413	21.7	1747.4	65	3	0
4.48E+09	0	18	5	0	3	8	2	1	1	0	14905	6.8	17974	233	52	-3282	-1647	23.3	1490.5	18	3	1
4.47E+09	1	22	7	1	9	4	10	1	1	0	18132	7	18076	228	40	2883	2111	22.8	1813.2	21	4	0
4.52E+09	1	17	4	1	5	2	6	2	1	1	16053	6.8	17406	218	51	1703	1696	21.8	1605.3	15	2	0
4.5E+09	0	12	3	0	6	4	7	0	0	0	16438	6.8	16989	203	32	608	45	20.3	1643.8	19	2	1
4.5E+09	0	15	2	1	5	4	4	0	0	0	13634	6.8	16355	201	45	71	-359	20.1	1363.4	14	1	0
4.51E+09	1	40	2	0	6	4	4	0	0	0	16085	7.4	19076	204	56	275	872	20.4	1608.5	15	1	1
4.52E+09	0	24	3	0	3	4	6	0	0	0	14981	6.8	17082	235	52	-809	638	23.5	1498.1	17	2	1
4.52E+09	0	12	4	0	2	8	2	1	1	0	15421	7.2	19276	217	84	-2771	-1125	21.7	1542.1	14	2	1

redFirstBl	redKills	redDeaths	redAssist	redEliteM	redDragor	redHerald	redTower	redTotalG	redAvgLe	redTotalE	redTotalM	redTotalU	redGoldD	redExperi	redCPerf	redGoldPe
0	6	9	8	0	0	0	0	16567	6.8	17047	197	55	-643	8	19.7	1656.7
1	5	5	2	2	1	1	1	17620	6.8	17438	240	52	2908	1173	24	1762
1	11	7	14	0	0	0	0	17285	6.8	17254	203	28	1172	1033	20.3	1728.5
1	5	4	10	0	0	0	0	16478	7	17961	235	47	1321	7	23.5	1647.8
1	6	6	7	1	1	0	0	17404	7	18313	225	67	1004	-230	22.5	1740.4
1	3	5	2	0	0	0	0	15201	7	18060	221	59	-698	-101	22.1	1520.1
0	6	7	9	0	0	0	0	14463	6.4	15404	164	35	-2411	-1563	16.4	1446.3
1	13	5	11	1	1	0	0	17920	6.6	16938	157	54	2615	800	15.7	1792
1	7	7	5	2	1	1	1	18380	7.2	19298	240	53	1979	771	24	1838
0	5	4	4	0	0	0	0	16605	6.8	18379	247	43	1548	1574	24.7	1660.5
0	4	4	5	1	1	0	0	15143	7.2	18196	216	51	-331	1585	21.6	1514.3
0	11	11	9	0	0	0	0	18200	7	19142	188	52	1505	635	18.8	1820
0	1	7	1	0	0	0	0	14591	6.8	17443	240	50	-3274	-1659	24	1459.1
1	9	4	11	0	0	0	0	18393	7.2	18354	229	51	3414	1141	22.9	1839.3
0	4	4	3	0	0	0	0	16192	7	18083	242	48	470	187	24.2	1619.2
1	7	3	13	0	0	0	0	17011	7.2	18778	237	51	1996	1804	23.7	1701.1
0	2	10	2	1	1	0	0	14505	6.8	17484	210	64	-5228	-3378	21	1450.5
1	7	3	6	0	0	0	0	16827	6.8	18233	218	53	1975	1345	21.8	1682.7
0	4	5	3	2	1	1	0	15400	6.6	16866	228	52	-882	-512	22.8	1540
1	9	5	9	0	0	0	0	18149	7.4	20697	219	56	3155	2773	21.9	1814.9
0	6	11	4	0	0	0	0	16080	7	17629	225	48	-2526	-1625	22.5	1608
1	13	4	14	1	1	0	0	18713	7.2	19766	209	45	2835	2821	20.9	1871.3
1	6	4	10	1	1	0	0	17027	7	18129	231	60	1254	567	23.1	1702.7
1	3	4	3	0	0	0	0	15500	7	17936	220	64	-406	-1283	22	1550
1	9	6	13	0	0	0	0	17535	6.6	16730	197	36	1020	-493	19.7	1753.5
1	8	4	13	1	1	0	0	17887	7	17114	221	36	2472	1067	22.1	1788.7
1	6	8	3	0	0	0	0	15918	6.8	17909	233	36	-1260	55	23.3	1591.8
0	4	7	2	1	1	0	0	15369	6.6	16967	202	49	-1857	-2015	20.2	1536.9
0	6	8	11	0	0	0	0	16154	6.8	17849	223	40	-1320	413	22.3	1615.4
1	8	3	12	0	0	0	0	18187	7.4	19621	225	52	3282	1647	22.5	1818.7
0	4	9	2	0	0	0	0	15249	6.2	15965	227	36	-2883	-2111	22.7	1524.9
0	2	5	2	0	0	0	0	14350	6.4	15710	206	48	-1703	-1696	20.6	1435
1	4	6	4	1	1	0	0	15830	6.8	17074	206	44	-608	85	20.6	1583
0	4	5	4	1	1	0	0	15563	6.8	16714	222	38	-71	359	22.2	1556.3
1	4	6	5	1	1	0	0	15810	7.2	18204	223	62	-275	-872	22.3	1581
1	4	3	7	2	1	1	0	15790	6.8	16444	218	50	809	-638	21.8	1579
1	8	2	5	0	0	0	0	18192	7.4	20401	234	64	2771	1125	23.4	1819.2

## 3 Specifications

### AdaBoost

Our first classifier was AdaBoost, AdaBoost is essentially a Random Forest classifier with some very key differences and methodologies. AdaBoost utilizes a forest of trees with only two leaves(stumps) and each stump has an amount of say determined by how correct its classification was. This then influences the next stumps amount of say by adjusting the weights assigned to the records and AdaBoost utilizes a weighted gini index to determine which feature will be split on to create the stump relative to the weights assigned. At first the weights are assigned as  $1/n$  and change through the course of the algorithms execution.

### Multi-Layer Perceptron

We also used a Mult-Layer Perceptron as one of our classifiers. A "MLP" is a type of neural network that utilizes multiple hidden layers with nodes in those layers and weights connected the nodes between layers. How this algorithm learns is that after completing a feed forward process the algorithm will then gather the error information and change the weights accordingly by means of back propagating through all layers and nodes again. Based on the solver type the weight changes will change weight after each phase of back propagation or another method. Doing this for the multiple layers and nodes allows the discovery of patterns in the data for more in depth predictions.

### Support Vector Machine

The final algorithm that we chose to use in our experimentation was a Support Vector Machine (SVM). An SVM is similar to a linear regression model, in that its objective is to find an optimal line of best fit (or optimal hyperplane in the case of 3-dimensional data or higher). A set of weights is used to calculate the weighted sum using the model input, which would produce the hyperplane. Afterwards, the error is calculated and the weights are adjusted accordingly using a method such as gradient descent. This will produce a hyperplane that separates the data, however, the optimal hyperplane (or the ultimate goal) is defined as the hyperplane that maximizes the amount of margin between the hyperplane and the first closest data point. Classification then becomes a simple task, input is fed into the model, and the predicted class becomes the corresponding side of the hyperplane where the output falls.

## 4 Implementation

### AdaBoost Implementation

The AdaBoost hyper-parameters that were chosen to adjust were the algorithm SAMME or SAMME.R, Learning Rate, and Number of Estimators. The in-

teraction of these hyper-parameters was also measured and experimented with. However the interaction between all 3 was too costly to get conclusive results, each run took an enormous amount of time to run. Overfitting and underfitting were monitored through accuracy and loss, if the accuracy stayed unmoving and the loss began to increase it was a sure sign of overfitting and vice versa for underfitting when going to low with the hyper-parameters and it wasn't fitting the data well at all.

## Multi-Layer Perceptron Implementation

For the MLP we established the default value as the a three hidden layers of 5 nodes, using an adam solver, a relu activation function and learning rate of 0.0001. Upon testing we altered the values for each of these described hyper-parameters. Our goal of this was to find if changing these values from their defaults would increase the accuracy. The answer to this question, right off the bat, is that it's a bit hard to tell. Our split method randomly splits the data so you never receive the same data every time. Also, with our data being so large and diversified a lot of variation could occur on how the algorithm takes in the inputs and learns from it. Both these reason causes each training session of the algorithm to produce wildly different results and this cause us to get an accuracy of around 74%, with slight variation, almost every time. Further testing would be needed, perhaps without allowing as much variation to occur, to allow accurate analysis.

## Support Vector Machine Implementation

The SVM hyperparameters that we selected were: C-regularization parameter, kernel function, the degree of the polynomial kernel, and maximum iterations. The aim was to determine which kernel function produced the best results and why. We found that during our experimentation, the linear or rbf kernel functions consistently produced the best accuracy values given our datasets. On the otherhand, the polynomial and sigmoid kernel functions performed the worst. This was an interesting observation to us, as we did not expect the linear and rbf kernels to outperform the polynomial function, since our data had a high number of features. As for the other hyperparameters, only slight changes in performance were observed whenever they were altered. Our conclusion was that the kernel function was the most important hyperparameter to change, which would result in the greatest change in model performance. Aside from a few other interesting observations in the patterns created by the other hyper-parameters as we tweaked them, the SVM produced very consistent accuracy values for both datasets over dozens of runs.

## 5 Conclusions

The Classifiers and their best accuracy for the best hyperparameters per dataset:

### League Dataset Accuraries

Classifier	Accuracy
AdaBoost+SAMME+300est	73.83%
MLP+40nodes	73.5%
SVM+1.9C+linear	72.7%

### Spam Dataset Accuracies

Classifier	Accuracy
AdaBoost+SAMME+500est	94.64%
MLP+5nodes	77.28%
SVM+4C	94.2%