

ABSTRACT

Five machine learning models were investigated along with a novel feature design to predict the outcomes and popular betting metrics of NBA basketball games. The most predictive feature design consisted of the statlines of the top 3 players of each team in their last 3-4 games. The support vector machines classifier performed the best in predicting game outcomes, achieving 62.6% accuracy on the test set. The logistic regression and quadratic discriminant analysis models were 57.8% and 56.4% accurate respectively. The exponential neural network (NN) score predictor achieved 59.1% accuracy and reproduced the score distribution of the test set well, indicating that it may be effective in predicting the over/under for NBA games. The softmax NN score predictor only achieved 57.6% accuracy but it managed to reproduce the margin of victory distribution well. Ultimately, a mismatch between the train set and the dev/test sets as well as an overall lack of data likely led to the unremarkable performance of the five models.

DATASET AND FEATURES

VISITOR: Los Angeles Lakers (20-3)

	POS	MIN	FG	FGA	3P	3PA	FT	FTA	OR	DR	TOT	A	PF	ST	TO	BS	+/-	PTS
23 LeBron James	F	33:45	11	23	4	9	5	7	1	6	7	8	1	0	3	1	21	31
3 Anthony Davis	F	32:09	12	21	2	6	13	15	0	9	9	2	2	2	3	3	6	39
7 JaVale McGee	C	15:25	6	7	0	0	1	1	1	2	3	0	3	0	0	2	5	13
14 Danny Green	G	21:08	1	3	1	3	0	0	0	5	5	2	2	1	1	1	5	3
1 Kentavious Caldwell-Pope	G	27:01	1	3	1	3	2	2	0	0	0	5	1	1	0	0	8	5
9 Rajon Rondo		14:32	2	2	2	2	0	0	0	3	3	3	0	2	2	1	13	6
0 Kyle Kuzma		24:26	6	13	3	6	0	0	0	6	6	1	2	0	2	0	19	15
30 Troy Daniels		12:37	3	5	3	5	0	0	0	1	1	1	3	1	1	0	7	9
4 Alex Caruso		23:46	2	4	1	1	3	4	0	1	1	3	4	1	1	0	13	8
39 Dwight Howard		21:20	1	1	0	0	3	4	1	9	10	0	5	2	2	0	14	5
2 Quinn Cook		10:07	1	3	0	1	0	0	0	0	0	1	0	0	1	0	0	2
10 Jared Dudley		03:44	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0

Figure 1: Example of a box score containing many of the stats utilized in this project.

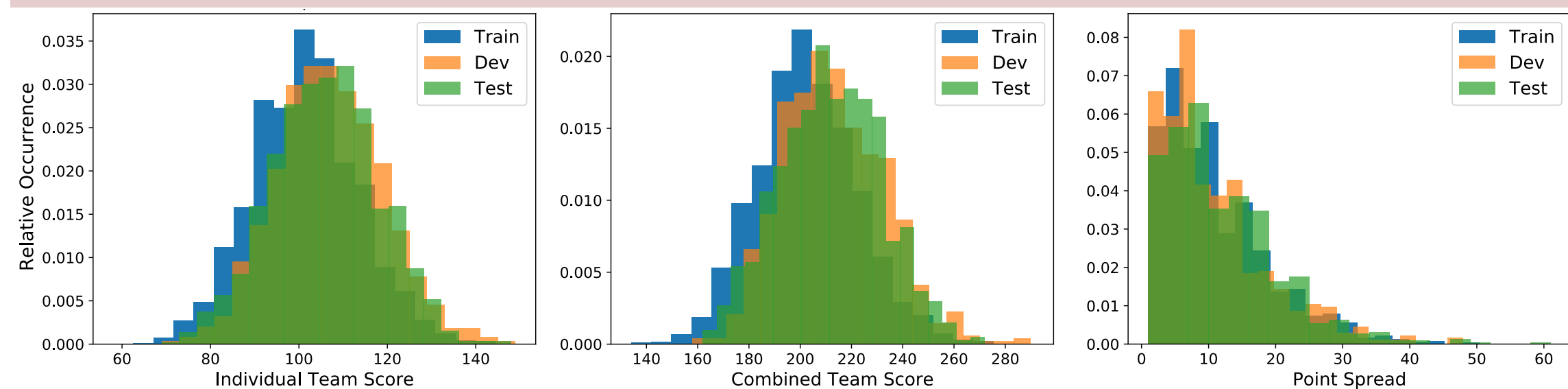


Figure 2: Histograms for individual team scores (left), combined scores (middle) and point spreads (right) from the NBA game corpus used in this project.

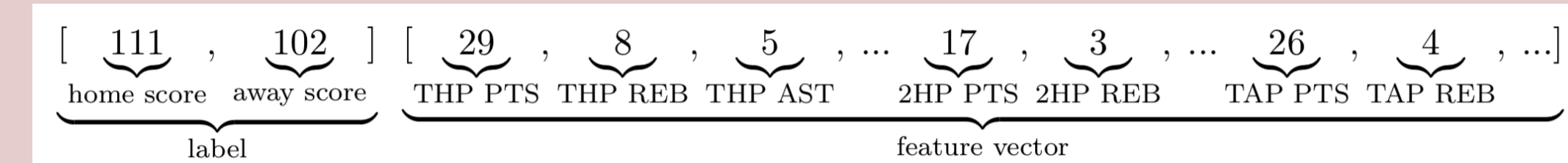


Figure 3: Example of a non-standardized feature vector for $n_g = 1$ and $n_p = 2$, where THP is top home player, 2HP is second-ranked home player, and TAP is top away player, and the stats present are the statlines for the players in question from the last game that they played

- The base dataset for this work is the complete statline for every player in each game from the 2012-2013 NBA season to the 2017-2018 season.
- A 80/10/10 split for train/dev/test was chosen, which results in dataset sizes of about ~5800/725/725 games depending on which feature parameters are chosen. The train set is the first 80% chunk of the games by chronological order and the test/dev set examples are sampled uniformly from the remaining 20%.
- Each feature vector contains the statlines of the top n_p players (ranked by Pts) that are on the roster for the game in question in their previous n_g games for each team.

MODELS

- In this work, five distinct machine learning models are explored.
- Three binary classifier models are employed via SKLearn; logistic regression, support vector machines, and gaussian discriminant analysis w/ different covariance matrices for each class, giving rise to quadratic decision boundaries.
- Two score prediction neural networks are implemented with Keras, one with an exponential output layer and the other with a softmax output layer.
- Due to the small size of the dataset, HPC resources were not required and all models were trained on a personal laptop

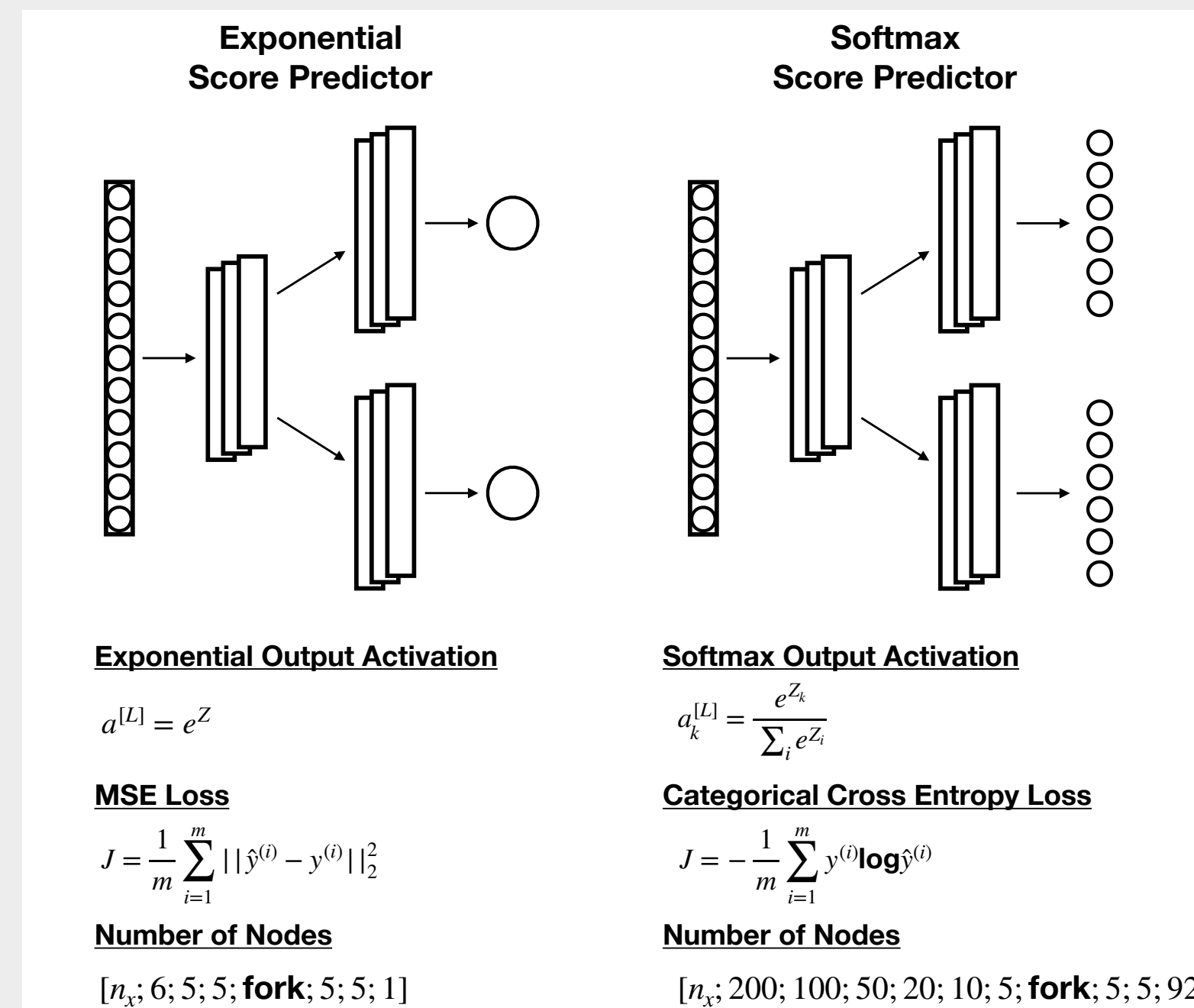


Figure 4: NN architectures for the two score predictor models tested in this work.

RESULTS

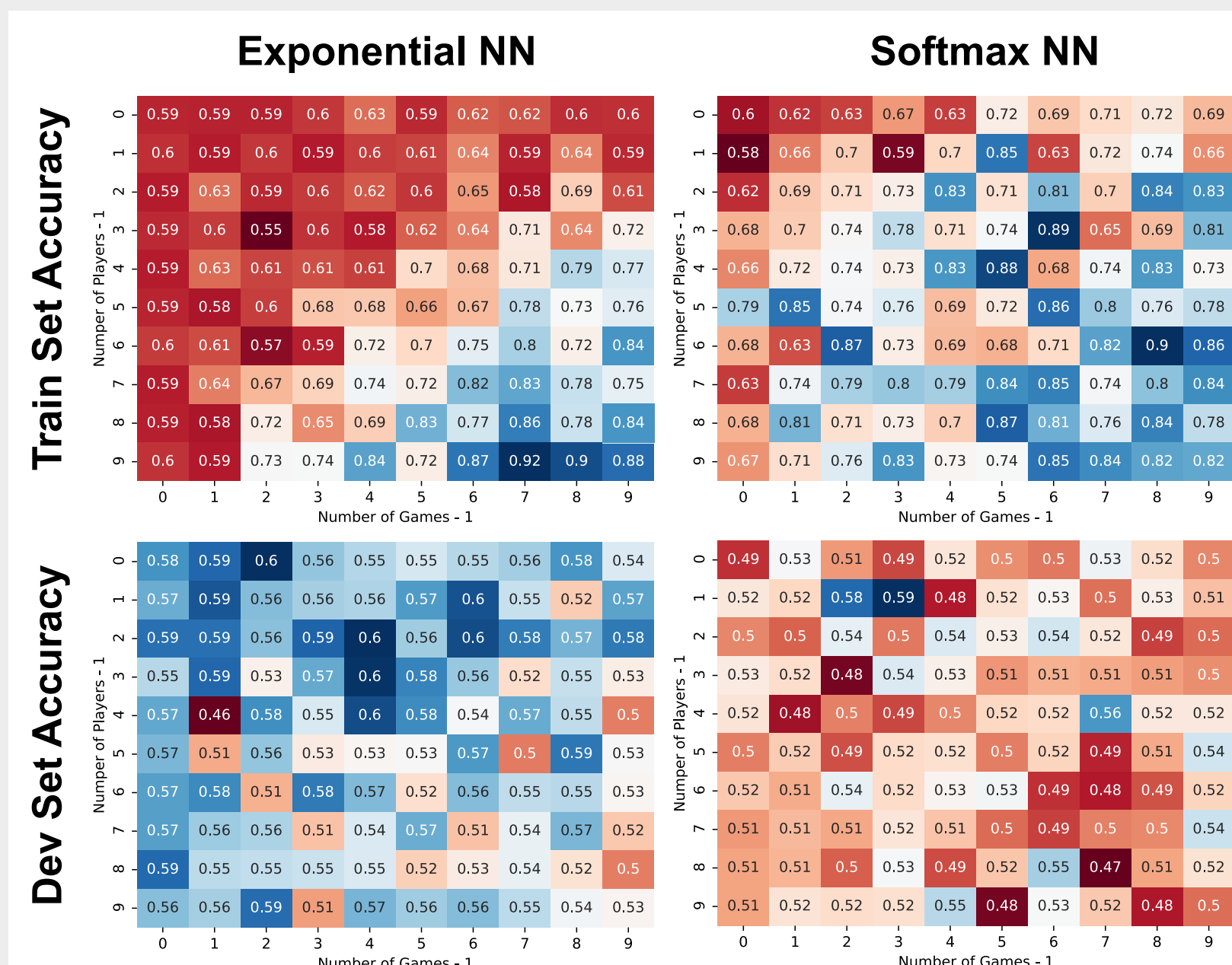


Figure 5: Dev set performance of the 2 score predictor models for each n_g, n_p pair.

RESULTS (CONT)

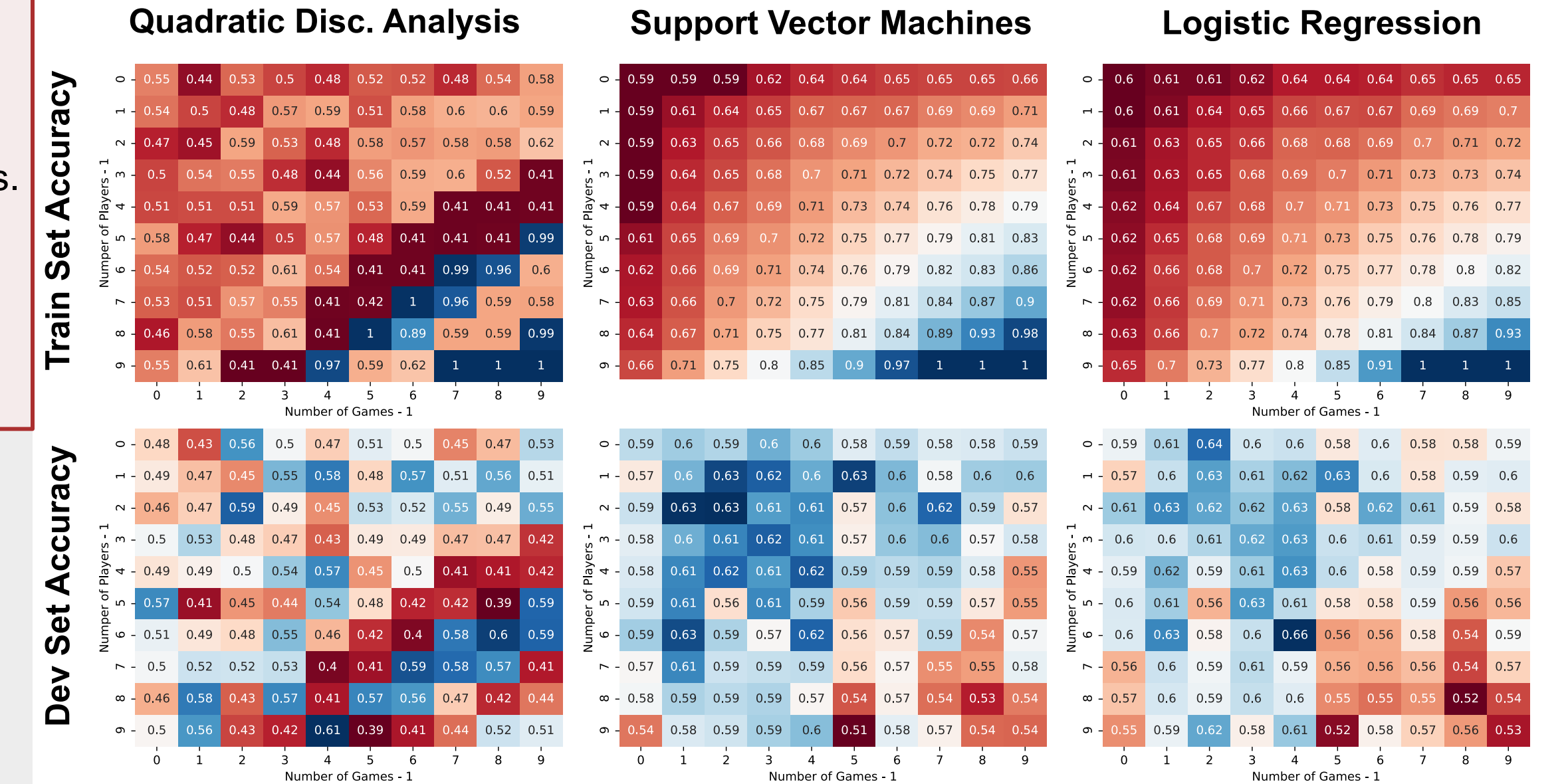


Figure 7: Dev set performance of the 3 binary classifier models for each n_g, n_p pair.

Model	n_g, n_p pair	Train Set Accuracy	Dev Set Accuracy	Test Set Accuracy
QDA	(4,10)	97%	61%	56.5%
SVM	(3,3)	65%	63%	62.6%
LR	(5,7)	72%	66%	57.8%
Exponential NN	(4,3)	62%	60%	59.1%
Softmax NN	(4,2)	59%	59%	57.6%

Figure 8: Accuracies on the train/dev/test sets for each model based on best n_g, n_p pair determined above.

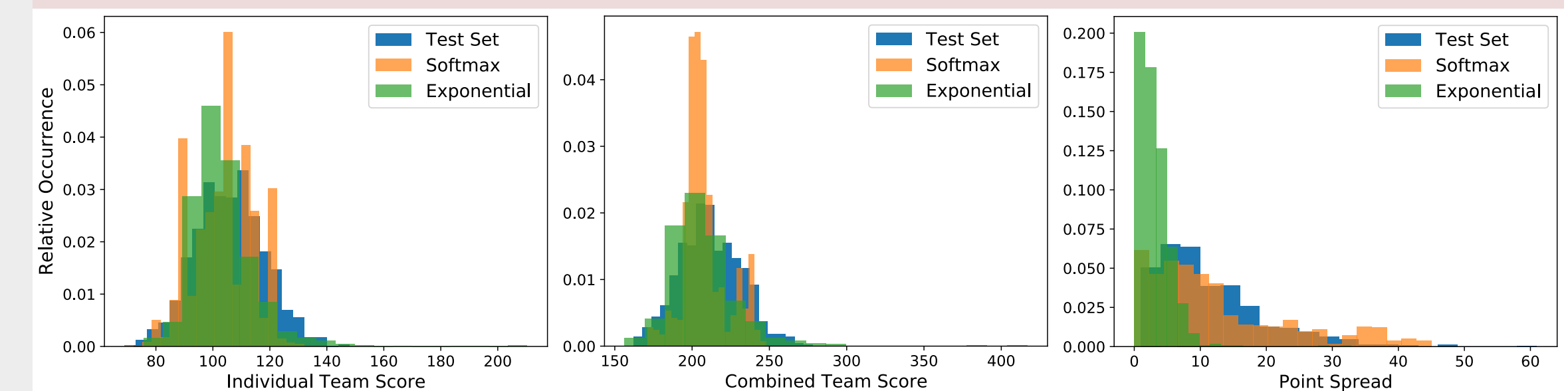


Figure 9: Histograms for individual team scores (left), combined scores (middle) and point spreads (right) from the true test set values and the predictions from each model.

CONCLUSIONS/FUTURE WORK

- The best features contained the stats of each team's top ~3 players in their last 4 games.
- The SVM classifier achieved a test set performance on par with human experts.
- The exponential score predictor reproduced the score distribution of the test set quite well, indicating that it may be effective in predicting the over/under for NBA games. The softmax score predictor performed relatively poorly on the previous two tasks but managed to reproduce the margin of victory distribution remarkably well.
- Omission of some stats along with the inclusion of other stats that weren't present in the dataset would likely improve performance by shrinking the size of the feature vectors and making the information present more potent.
- Ultimately, the main issue is the lack of training data and the mismatch between the train and dev/test sets. To tackle both of these problems, it may be effective to artificially generate data via NBA basketball video games.