

Fourth Down Decision Predictor

Anthony Clavette & Lexington Johnson

Abstract:

The goal of this project is to correctly predict when a team will go for it on fourth down. When the defense is prepared and stops a team on a fourth down attempt, it gives the ball to their offense which helps them win. We want to help defenses have a continuous awareness of the opposing team's probability of going for it on fourth down or kicking the ball away. Despite being just a game, professional sports results can have major impacts on people's jobs and communities, so the risky and poor decisions made by defensive coaches on fourth down must be nearly perfect. We found that the score differential, the yards to go for a first down, the time left in the game, and the team's position on the field most heavily contribute to their fourth down decision of kicking or going for it.

Introduction:

As technology becomes exponentially more advanced, we have seen a significant rise in innovation with offensive schemes than just ten years ago. Since the defense is always one step behind the offense, more time and algorithmic decision-making is required for defensive coordinators to prepare their side. Our model can provide defensive coaches with information that will help them be more aware of what contributes to the opposing team's fourth down decision. Professional football is a billion-dollar entertainment industry that floods United States economics, media, and pop culture. Game outcomes have been proven to be capable of causing major uproar in cities and on the internet, so these split-second decisions that coaches have to make can have an enormous impact on the NFL community. Typically, other works using this dataset or analyzing similar data focus on enhancing the offensive side of decision making throughout the game. On the other hand, our work focuses on preparing the defense and giving them more time for smaller details in fourth down scenarios. We can do this by determining what features of the game are the most important and revealing hidden situational patterns that are invisible to the naked eye.

Background:

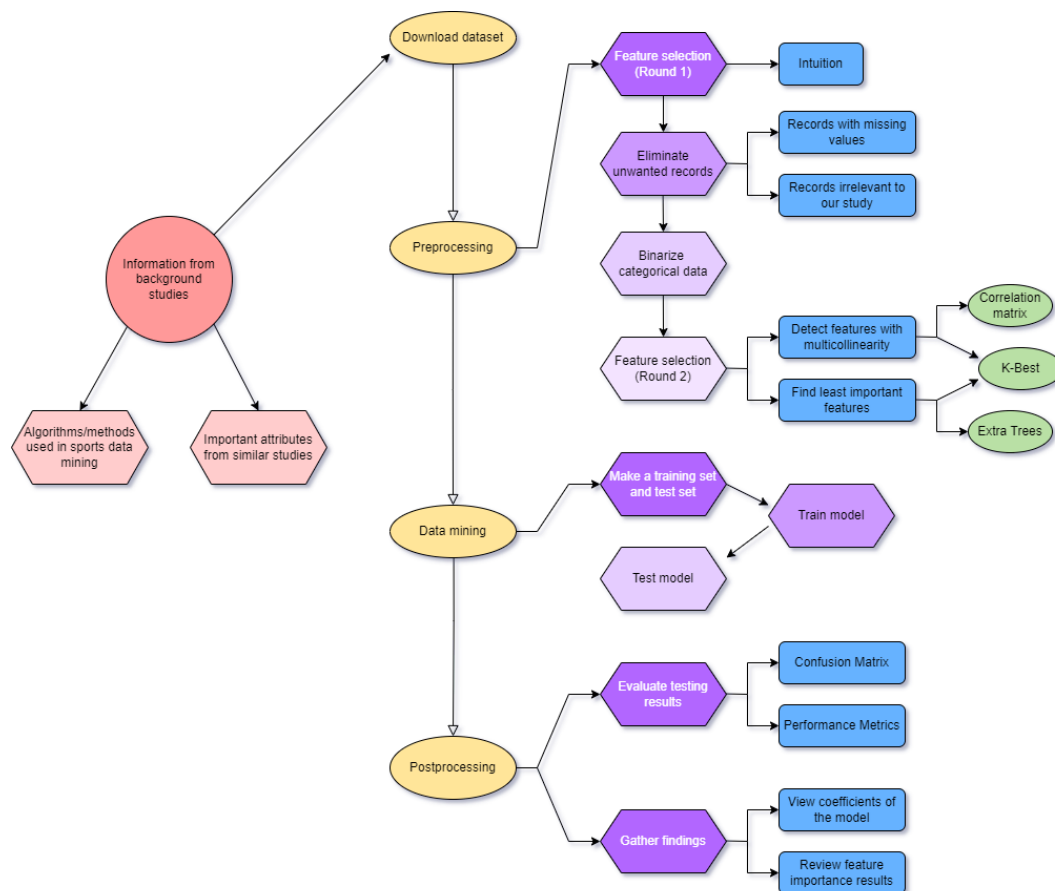
The sports data industry has stretched to expensive heights and made incredible advancements, which means many studies have solved problems similar to ours. Studies helpful to our progress include statistical studies on American football strategy, fourth down strategies, and the most important attributes to focus our model around. These studies also helped us adopt a plan and healthy habits for the project, such as how to properly prepare the dataset for usage in preprocessing algorithms and classifiers. These studies also gave us examples of how to analyze the results of a predictive model and what are the most meaningful performance metrics to a problem like ours.

To summarize a few background studies on fourth down decisions, their models all shared the following features: yards to first down, yard line on the field, and time remaining. Therefore, we were able to assume that these would have the largest weight in our model, and we

should pay the most attention to them. Some of them focused on some more obscure elements to the problem, such as analyzing how fourth down decisions and success changes as the season goes on. Aside from predictive models, another useful study mapped recommended fourth down decisions according to success probability and the attributes mentioned above. This allowed us to incorporate a feature describing the team's win probability because this study outlined its correlation with fourth down success. Overall, reading studies of similar datasets or similar problems being solved gave us an advantage with understanding the importance of our dataset's features and gave useful insight on effective methodology and work progress.

Methodology:

The framework for our research and processes is shown in the flowchart below. After reading background studies and downloading the dataset, our work had three distinct steps: preprocessing, data mining, and post processing.

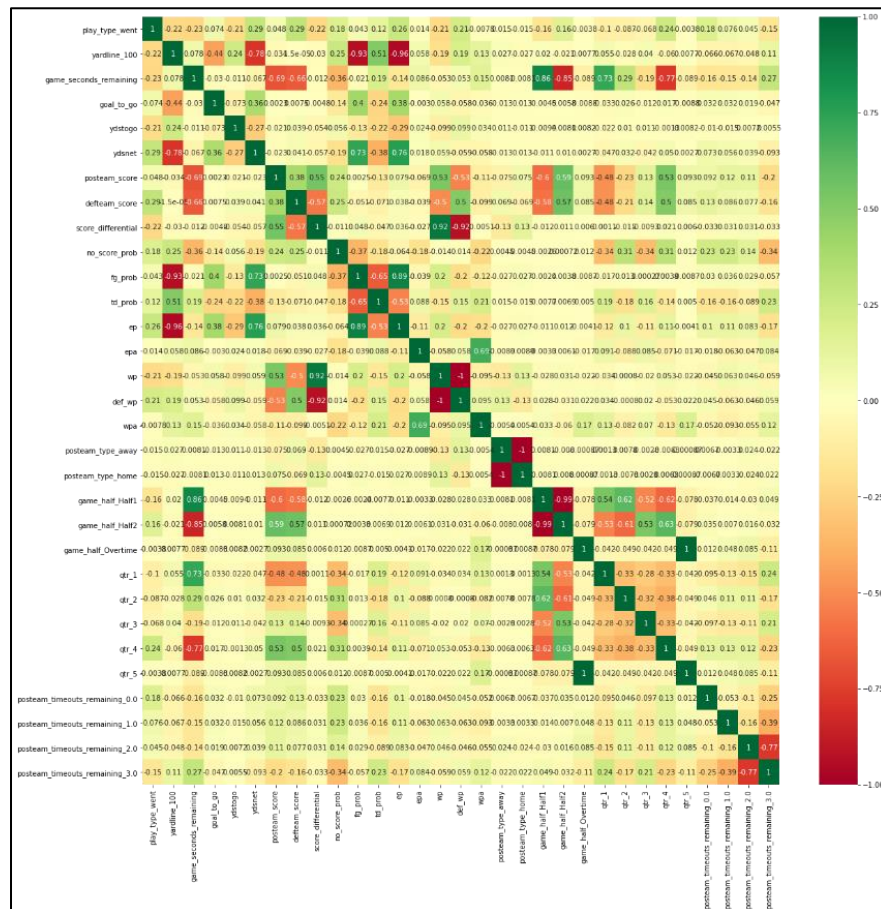


During preprocessing, we removed over 200 attributes by intuition and then needed to use algorithms to be able to keep reducing dimensionality. Supported by Pandas and Scikit-learn, we used a correlation matrix, a K-Best classifier, and an Extra Trees classifier to reveal the least important features and any multicollinearity that may exist. To begin the data mining step, we used built-in methods from Scikit-learn to build our training and test sets, build a logistic regression model, train the model, and test the model on the test data. In order to meaningfully

interpret our trained logistic regression model, we wanted to view the results of our testing and information we found about our attributes. To evaluate the results of our model, we used algorithms from Scikit-learn to view the confusion matrix, performance metrics, and ROC/AUC. Lastly, we looked at the model's weight coefficient for each attribute and we reviewed the results of the feature importance algorithms from preprocessing.

Data:

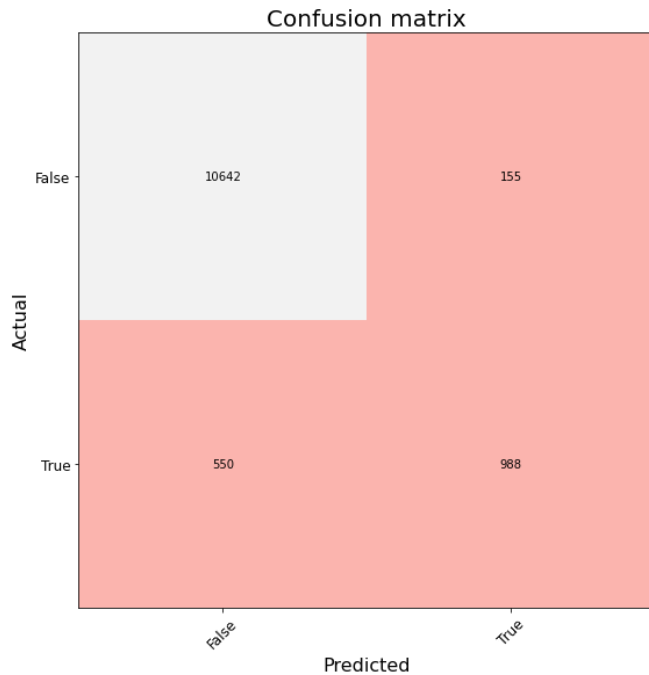
The dataset was originally a 700 MB collection of every NFL play that happened from 2009-2018. There were 255 attributes describing the nature of all 450,000+ records in the dataset, which would be a lot of data for us to efficiently process. Fortunately, many attributes and records were completely irrelevant to our study and could be removed, we just had to find them. We started with feature selection, removing the attributes that are obviously unnecessary and do not contribute to the independent variable we were targeting. We were able to remove about 200 attributes this way. To clean the data more before running it through any preprocessing algorithms, we removed all records with missing values and also irrelevant records (unfinished plays, penalties, or plays that happened on first, second, and third down). To go deeper into feature selection, we produced the following correlation matrix to find any interesting correlations or multicollinearity among attributes:



		month predicted position																		
		37	48	42	65	78	87	91	104	113	126	185	212	239	223	239				
369	276	298	315	328	323	352	362	375	402	435	458	418	418	458	466	476				
400	534	515	544	558	574	584	608	613	631	635	657	658	663	662	662	662				
401	534	515	544	558	574	584	608	613	631	635	657	658	663	662	662	662				
781	708	783	792	797	801	811	834	858	868	874	878	883	885	894	932	932				
801	708	783	792	797	801	811	834	858	868	874	878	883	885	894	932	932				
811	1112	1122	1143	1153	1153	1192	1197	1208	1209	1213	1221	1224	1228	1230	1232	1232				
1208	1277	1285	1288	1294	1298	1308	1309	1312	1344	1347	1352	1352	1358	1359	1360	1360				
1277	1278	1271	1271	1271	1271	1271	1271	1271	1271	1271	1271	1271	1271	1271	1271	1271				
1403	1488	1489	1471	1473	1476	1473	1473	1473	1473	1473	1473	1473	1473	1473	1473	1473				
1403	1488	1489	1471	1473	1476	1473	1473	1473	1473	1473	1473	1473	1473	1473	1473	1473				
1602	1606	1607	1606	1606	1606	1606	1606	1606	1606	1606	1606	1606	1606	1606	1606	1606				
1606	1608	1602	1620	1620	1620	1620	1620	1620	1620	1620	1620	1620	1620	1620	1620	1620				
1606	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607	1607				
1074	1077	1091	1083	1083	1083	1083	1083	1083	1083	1083	1083	1083	1083	1083	1083	1083				
1218	2218	2226	2229	2236	2236	2236	2236	2236	2236	2236	2236	2236	2236	2236	2236	2236				
2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403	2403				
2501	2508	2502	2540	2543	2605	2672	2688	2684	2684	2684	2684	2684	2684	2684	2684	2684				
2501	2508	2502	2540	2543	2605	2672	2688	2684	2684	2684	2684	2684	2684	2684	2684	2684				
2508	2509	2548	2584	2584	2584	2584	2584	2584	2584	2584	2584	2584	2584	2584	2584	2584				
3078	3071	3076	3084	3083	3110	3112	3112	3112	3112	3112	3112	3112	3112	3112	3112	3112				
3078	3071	3076	3084	3083	3110	3112	3112	3112	3112	3112	3112	3112	3112	3112	3112	3112				
3132	3141	3142	3158	3152	3154	3154	3154	3154	3154	3154	3154	3154	3154	3154	3154	3154				
3132	3141	3142	3158	3152	3154	3154	3154	3154	3154	3154	3154	3154	3154	3154	3154	3154				
3573	3574	3602	3610	3628	3646	3678	3680	3687	3697	3702	3719	3732	3737	3737	3737	3737				
3573	3574	3602	3610	3628	3646	3678	3680	3687	3697	3702	3719	3732	3737	3737	3737	3737				
3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746	3746				
3958	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978	3978				
4129	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438				
4129	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438				
4129	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438	4438				
4608	4605	4527	4513	4514	4562	4567	4568	4576	4582	4595	4608	4618	4638	4648	4664	4664				
4608	4605	4527	4513	4514	4562	4567	4568	4576	4582	4595	4608	4618	4638	4648	4664	4664				
4634	4687	4851	4858	4858	4879	4877	4883	4881	4881	4885	4895	4907	4917	4933	4944	4944				
4634	4687	4851	4858	4858	4879	4877	4883	4881	4881	4885	4895	4907	4917	4933	4944	4944				
4905	4973	4905	4990	5004	5008	5008	5022	5011	5015	5017	5023	5032	5048	5068	5087	5097				
4905	4973	4905	4990	5004	5008	5008	5022	5011	5015	5017	5023	5032	5048	5068	5087	5097				
5087	5087	5087	5112	5112	5112	5112	5112	5112	5112	5112	5112	5112	5112	5112	5112	5112				
5208	5322	5248	5338	5341	5266	5273	5278	5288	5288	5292	5305	5311	5315	5321	5321	5321				
5208	5322	5248	5338	5341	5266	5273	5278	5288	5288	5292	5305	5311	5315	5321	5321	5321				
5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307	5307				
5402	5407	5482	5510	5528	5528	5552	5552	5555	5557	5557	5568	5578	5585	5595	5608	5608				
5402	5407	5482	5510	5528	5528	5552	5552	5555	5557	5557	5568	5578	5585	5595	5608	5608				
5611	5621	5643	5655	5661	5663	5679	5706	5708	5719	5757	5759	5781	5803	5810	5813	5813				
5611	5621	5643	5655	5661	5663	5679	5706	5708	5719	5757	5759	5781	5803	5810	5813	5813				
5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813				
5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813	5813				
6308	6305	6054	6064	5976	5995	6016	6016	6024	6024	6024	6024	6024	6024	6024	6024	6024				
6308	6305	6054	6064	5976	5995	6016	6016	6024	6024	6024	6024	6024	6024	6024	6024	6024				
6308	6305	6054	6064	5976	5995	6016	6016	6024	6024	6024	6024	6024	6024	6024	6024	6024				
6328	6325	6323	6356	6368	6368	6388	6388	6393	6402	6402	6406	6412	6414	6428	6438	6448				
6328	6325	6323	6356	6368	6368	6388	6388	6393	6402	6402	6406	6412	6414	6428	6438	6448				
6407	6403	6407	6408	6409	6407	6412	6419	6428	6454	6454	6454	6463	6463	6463	6463	6463				
6407	6403	6407	6408	6409	6407	6412	6419	6428	6454	6454	6454	6463	6463	6463	6463	6463				
6603	6613	6623	6626	6648	6688	6687	6693	6693	6693	6693	6693	6693	6693	6693	6693	6693				
6603	6613	6623	6626	6648	6688	6687	6693	6693	6693	6693	6693	6693	6693	6693	6693	6693				
6603	6613	6623	6626	6648	6688	6687	6693	6693	6693	6693	6693	6693	6693	6693	6693	6693				
7078	7086	7086	7086	7086	7102	7112	7117	7138	7149	7170	7180	7190	7190	7190	7190	7190				
7078	7086	7086	7086	7086	7102	7112	7117	7138	7149	7170	7180	7190	7190	7190	7190	7190				
7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246	7246				
7405	7409	7409	7409	7442	7458	7463	7467	7496	7496	7496	7503	7512	7513	7514	7554	7554				
7405	7409	7409	7409	7442	7458	7463	7467	7496	7496	7496	7503	7512	7513	7514	7554	7554				
7636	7636	7636	7636	7636	7644	7644	7644	7644	7644	7644	7644	7644	7644	7644	7644	7644				
7636	7636	7636	7636	7636	7644	7644	7644	7644	7644	7644	7644	7644	7644	7644	7644	7644				
7806	7808	7915	7920	7912	7935	7948	7985	8003	8013	8024	8034	8036	8041	8042	8047	8047				
7806	7808	7915	7920	7912	7935	7948	7985	8003	8013	8024	8034	8036	8041	8042	8047	8047				
8254	8258	8257	8257	8258	8258	8258	8257	8312	8312	8312	8312	8312	8312	8312	8312	8312				
8254	8258	8257	8257	8258	8258	8258	8257	8312	8312	8312	8312	8312	8312	8312	8312	8312				
8731	8718	8821	8821	8821	8821	8821	8827	8879	8873	8883	8884	8951	8958	8951	8951	8951				
8731	8718	8821	8821	8821	8821	8821	8827	8879	8873	8883	8884	8951	8958	8951	8951	8951				
8731	8718	8821	8821	8821	8821	8821	8827	8879	8873	8883	8884	8951	8958	8951	8951	8951				
8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958				
8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958	8958				
9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064				
9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064				
9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064	9064				
9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156				
9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156				
9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156	9156				
9751	9756	9756	9756	9756	9756	9756	9756													

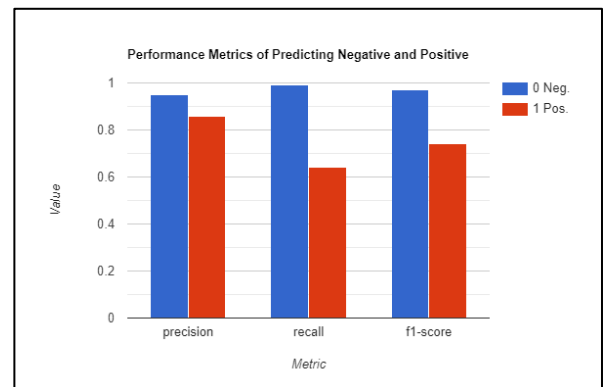
Results:

The model predicted that the team would attempt a play on fourth down (positive outcome) for 1143 instances of the test set's 12,344 total records, with 94% accuracy. The confusion matrix and performance metrics are shown below:



```
from sklearn.metrics import classification_report
print(classification_report(Y_test, predictions))
```

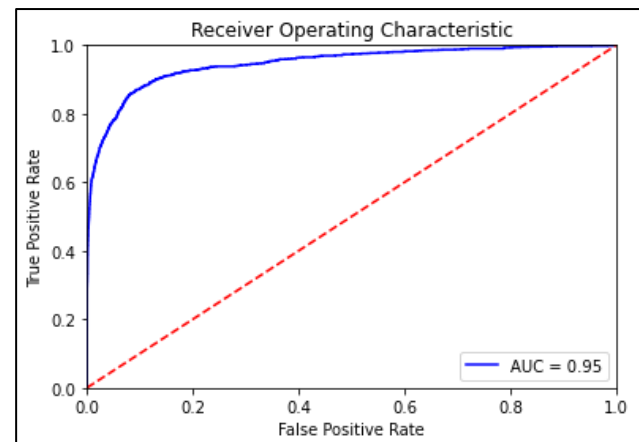
	precision	recall	f1-score	support
0	0.95	0.99	0.97	10797
1	0.86	0.64	0.74	1538
accuracy			0.94	12335
macro avg	0.91	0.81	0.85	12335
weighted avg	0.94	0.94	0.94	12335



These results show that the model clearly proved its competence, but they also show weak spots in it as well. Looking at the precision, recall, and f1-score for the two outcome predictions, we can see how much better the model was for predicting negative outcomes than positive outcomes. This means the high accuracy of the model (0.94) can be linked to the dataset having a vast majority of true negatives and the model's ability to correctly predict them.

The graph on the right is the ROC curve, which shows the TPR plotted against the FPR. This curve visualizes the model's discriminative ability, and the area under the curve is used as a measure for the trait. The area under our curve was 0.95, which shows a consistent proficiency in distinguishing between positive and negative outcomes.

As a result of these findings, we can conclude that although our model is better at correctly predicting negative outcomes than positive outcomes, it is still very good at distinguishing between both classes.



Conclusion:

We had some interesting results that came from data exploration and testing our model, but nothing was too shocking. We had a few ideas of what to expect in our findings because our background studies gave us an idea of important and unimportant features and we are also just familiar with American football.

Early on, an obstacle we faced was using Weka with our dataset and trying to familiarize ourselves with the software. We were intrigued by Weka's organized GUI, built-in libraries and algorithms, model summary output, and especially the visualization methods. Overall, the software seemed optimal for taking this project to a higher level with every resource we could ever need. The obstacle was format issues of some attribute values in our dataset. We spent hours troubleshooting by trying different methods of discretizing, binarizing, and other forms of preprocessing our data because Weka refused to build a logistic regression model using our dataset. After about a week of research and trying different solutions, we started using Pandas and Scikit-learn with Python instead; we had our data cleaned and our first logistic regression model built within just one long day of working on it. On the other hand, there were also challenges we were not able to solve but have the potential to extend and improve this project further. The biggest challenge is transforming "irrelevant" records into attributes. Specifically, instead of just removing all plays that happened on first/second/third down, we would only remove the ones that did not eventually reach fourth down. For each fourth down being examined, our goal was to store the results of the preceding first/second/third downs as another feature. We believe that the individual outcomes of first/second/third down are relevant to the outcome of their corresponding fourth down, so we look forward to figuring out the solution to this challenge in the near future.