



Basketbal	ll was inver	nted in 1891	by Dr. Jam	nes Naismith	h. At first,
the game	was very si	imple. The l	baskets wer	e made from	n peach
baskets, a	nd the ball	in use was	actually a so	occer ball. (	Over 100
years late	r, the game	grew to be	come much	more comp	lex with
innovatio	ns such as t	the three-po	oint line, a re	eplay system	n, and of
course, co	omprehensi	ve statistics	•		
Kobe Bry	ant is curre	ently one of	the top scor	rer's in the	NBA, and
one of the	e greatest of	f all time. H	le is the you	ingest playe	er to score
31,000 pc	oints, is curi	rently 4 <sup>th</sup> in	all-time sco	oring, and h	as won two
scoring ti	tles. What c	lefines Kob	e Bryant is	his undenia	ble ability
to put the	ball in the	hoop.			
Therefore	e, one quest	ion we wan	t to investig	gate is how 2	Kobe's prior
pertormat	nces attect	his next cor	ing perform	nance. In or	der to
answer th	at question,	, we will pro-	edict Kobe	Bryant's sco	oring
technique		liext game i	le plays by	using mach	me learning
teeninque	<b>.</b>				
		<u>EXPLC</u>	<u>DRATOF</u>	<u>RY</u>	
We first for that PTS is played, field	und the correlated v d goals, fiel	relation bet with parame ld goals atte	ween predic eters such as empted, ave	ctors and po s game start rage points	oints. We four ter, minutes from past 5
We first for that PTS is played, field games, and Game Starter	und the correlated v d goals, fiel GmSc (a m Minutes Played	relation bet with parame ld goals atteneasure of k Field Goals	ween predic eters such as empted, ave Cobe's produ Field Goal Attempts	etors and po s game start rage points uctivity dur Average Points	oints. We four ter, minutes from past 5 ing a game). Game Score
We first for that PTS is played, fiel games, and Game Starter 0.462004	und the correlated vo d goals, fiel GmSc (a m Minutes Played 0.485752	relation bet with parame ld goals attente neasure of k Field Goals 0.491646	ween predic eters such as empted, ave Cobe's produ Field Goal Attempts 0.519255	ctors and pors s game start rage points uctivity dur: Average Points 0.509064	oints. We four ter, minutes from past 5 ing a game). Game Score 0.456355
We first for that PTS is played, fiel games, and Game Starter 0.462004 nce we had tween our p gher than 0.3 is large num mbinations cided to exp alt with pos	und the correlated version of several version of se	relation bet with parame ld goals attented neasure of k Field Goals 0.491646 0.491646 irs of varial ot even include variables with models that earity.	ween predic eters such as empted, ave tobe's produ Field Goal Attempts 0.519255 anted to che pairs of va bles with co lude colline th another v t performed	etors and por s game start rage points uctivity dur: Average Points 0.509064 eck the corre- riables with orrelation his arity betwee variable. The variable elia	oints. We four ter, minutes from past 5 ing a game). Game Score 0.456355 elation value correlation gher than 0.7 en linear us, we imination and

INTRODUCTION

The above plot shows the distribution of the number of points Kobe scored over time. Larger game numbers indicate later games.

## DATA

We gathered seasonal data from the basketball-reference website. We built table for Kobe's personal performance data from years 1996 to 2012. Since we wanted to predict Kobe's performance for his next game, the opponent he is playing with is an important determining factor for his performance. Thus, we built tables of offensive and defensive data for Kobe's opponent team from years 1996 to 2012. We then merged Kobe's personal data with his opponent data by combining on the opponent team's name.

www.PosterPresentations.com

# Predicting Kobe's Performance Stat 154 Modern Statistical Prediction and Machine Learning Cheng Dong, Min Jo, Frank Zhang

Since the parameters describing each game is unknown prior to the game, we used a 5 game grouping scheme to create predictor variables. We took the mean of the statistics from the previous five games of the game we want to predict and used those as predictor variables. At this point, we have 1064 games and 57 variables.

We had a few problems within our data that made certain entries unusable for model building. We had to use regex to convert some of the parameters into workable formats. These parameters included age, time, home or away games, win/loss, etc.

Some other problems we experienced included dealing with NAs and Infs in the Kobe's personal data. These resulted in percentage predictors. For instance, field goal percentage is equal to field goal successes over field goal attempts and will yield in NA if field goal attempts is zero. Since the percentage predictors were already represented by the variables that yielded in the percentage, we decided to remove these columns due to redundancy and obvious collinearity. Our final data set consisted of 1048 games and 51 variables.

## **METHODS**

Cross validation helped us evaluate the different models we used. Since we have a set amount of data, it was very useful for us to be able to use different training sets to see how well our model really performed.

OLS & GLS						
Method	CV Error	Training	Testing			
OLS (Forward Selection)	6.628941	6.943838	6.74374			
OLS (Backward Selection)	6.490223	6.970129	6.6932			
GLS (Forward Selection)	6.950420	6.797303	6.805020			

In order to predict Kobe's points in the next game, we utilized forward/ backward stepwise method to reduce dimensions of the data by selecting predictor variables. Using the OLS method with these variables, we actually got the lowest error of the all methods which is 6.49. The error is the mean of absolute values of the differences. For GLS, the error was 6.95. It should be noted that we used variable selection techniques to prevent overfitting.



The above plot shows the distribution of the number of points Kobe scored over time. Larger game numbers indicate later games.

10-14 15-19 20-25 25-30

With KNN, we tried several K's, and found that as k increased, the overall error rate went down. However, after further investigation, the reason why the error rate went down was because the algorithm predicted more and more 6's. The number of 6's correctly predicted went up at a faster rate than the rate for the number of 1-5's correctly predicted went down. Therefore, the error rate decreased even though the predicted values shifted towards more 6's.

We decided that the best k for this dataset is actually 5 because while the error went down as k increased, it was for the wrong reason. The algorithm almost seemingly blindly predicted more and more 6's. **Benchmarks** One benchmark that we used was simply looking at the previous season's

scoring average as our prediction for the next game. Therefore, we could not predict any game for the first season and our predictions are the same for all games in a season. The average absolute value of the residual was 8.485.

Another benchmark that we used was using his running career average as the prediction. This means that our prediction was very volatile early in his career and becomes stable towards the end of his career. The average absolute value of the residual was 8.191 which meant that it id better than the other benchmark, but still worse than our machine learning techniques.

<u>PCR</u>						
К	CV Error	Error (standardized)	Training	Training (stand)	Testing	Testing (stand)
1	6.5736	7.180537		7.16173		7.0078
5	6.5636	7.177982		7.14285		6.99745
10	6.4981	7.126609		7.06924		6.96534
40	6.4023	7.128349		6.89359		6.90954

The objective of our principal components analysis was to find a linear transformation of a set of our 50 predictor variables into a new set denoted by P, where the new set has certain desirable properties.

1) The elements of P are uncorrelated with each other in the sample 2) Each element of P accounts for as much of the combined variance of the x's as possible, consistent with being orthogonal to the preceding p's. As a result, we noticed that, as n increases, the magnitude of differences (errors) decreases.



For Stag

LAR (Lea Regre



While ridge regression is used to deal with high dimensional data, this method smoothes the coefficients of the model out instead of setting them to zero and completely eliminating its effects. The pros of this model is that each parameter will always have some, no matter how small, effect on the response variable. This pro is also a con since this makes it impossible to ever completely rid of collinearity. Overall, ridge regression performed the best among the LARs-like models. The residuals are randomly distributed about zero.

In order to predict Kobe's scoring performance for the next game, we tried multiple methods. It is difficult to distinctively point out the best method because most methods yielded similar results. However, after completing this project, there are several improvements that could have affected our results.

-More predictor variables (Kobe's one-on-one defender, his past performances specifically against that team. Etc....)

However, based on our predictions, we can see that most algorithms predict that the average points Kobe will score is around 26 points, which is consistent with his career average. The average residual is around 6.8 which is not terrible and considerably better than any benchmarks.

Method		CV Error	CV Erroi (standardiz	ed)	
La	asso	6.570095	7.280089	9	
Ridge		6.178731	7.074153	3	
Forward Stagewise		7.701985	NA		
LAR (Le Regr	east Angle ession)	6.725991	7.326344	4	
thod	Training	Training (stand)	Testing	Testing (stand)	
SSO	6.72508	6.734532	7.515963	7.177202	
dge	6.785631	7.060684	6.819793	6.873079	
ward ewise	7.701985	6.785631	7.060684	6.873079	
ast Angle ession)	6.716766	6.732922	7.16425	6.968429	
		Residuals for Ridge			

50	100	150	200	250	300	
Index						

## Conclusion

-Apply our techniques to more players in order to better understand the accuracy of our methods

## Acknowledgements

Johann Gagnon-Bartsch, Dat Duong, Kobe "Black Mamba" Bryant