

A Comparison of Statistical Learning Methods for Regression

OLS, LASSO, bagging: A Comparison

Andrew Bates

December 20, 2018

Executive Summary

This is the executive summary.

1 Introduction

In 2001 Leo Breiman described two approaches to analyzing data with statistical models (Breiman 2001). In *data modeling* we specify a stochastic model for the data, one that has known theoretical properties, and the main purpose is usually to make inferences on the population of interest. In the other approach, what Breiman calls *algorithmic modeling*, the focus is less on the the structure of the data itself and more about the output of the model. We are not concerned about finding a model that satisfies the theoretical assumptions needed to make inferences. We are interested in whether the model can make accurate predictions on newly collected data. Data modeling is the method typically taught in statistics and is probably what most who have studied statistics think of when they think about modeling. However, there has been increased interest in algorithmic modeling recently¹ with some, including Breiman, diminishing traditional statistical modeling.

In this paper we compare Breiman's two modeling paradigms by analyzing Major League Baseball data with the goal of developing a model to predict a players salary. We examine one data model (linear regression), one algorithmic model (random forest), and one model at the intersection of the two approaches (LASSO). For each model, we discuss some advantages and disadvantages in terms of both predictive capability and interpretability. The primary aim is to construct a predictive model but, although prediction and interpretability are often seen at odds with one another (Breiman 2001, 206), in some situations one may be interested in finding a balance between the two.

2 Methods

In this analysis we use the `Hitters` data from the R package `ISLR` (James et al. 2017), a companion package to *An Introduction to Statistical Learning with Applications in R* (James et al. 2013) containing the data sets used in the book. The `Hitters` data contains information on 322 players from the 1986 and 1987 Major League Baseball (MLB) seasons. There are 19 covariates included in this data set that can mostly be broken down into two categories: performance metrics for the 1986 season (number of at bats, number of home runs, etc.), and performance metrics based on a given players career (career runs, career hits, etc.). There are 16 continuous variables and three

¹See <https://trends.google.com/trends/explore?date=all&geo=US&q=machine%20learning> for example which shows web interest in machine learning since 2004 (accessed 12/9/2018).

categorical variables. For the 1987 season we have the player's salary on opening day along with their league (American or National) at the beginning of the season.

The salary variable has 59 missing observations, 18% of the data. This was too many observations to ignore so we imputed the values using k-nearest neighbors before proceeding with the analysis. After examining histograms of the continuous variables, it was evident that transformations were in order. Salary, along with several covariates, were heavily right-skewed. We chose log transformations for these variables because it is a common technique and allows us to readily interpret linear regression coefficients. In all, 11 of the 20 variables were log transformed. All numeric variables were then subsequently centered about the mean and scaled by the standard deviation.

Prior to model fitting the data was split into a training and testing set with 20% reserved for testing. All three models were trained using 5-fold cross-validation with the final model chosen to be the one with the lowest root mean squared error (RMSE). In each cross-validation run for linear regression a stepwise procedure was used with Akaike Information Criteria (AIC) as the model selection criterion. Diagnostics were run on the model chosen via cross-validation and some covariates were subsequently omitted based on variance inflation factors and correlations between the covariates. Grid search was used for hyperparameter tuning of the lasso and random forest with 10 values in each grid. For the lasso the hyperparameter is the lasso penalty and for random forest the hyperparameter is the number of randomly selected covariates considered at each split. A comparison of predictive ability for the three models was made through RMSE on the testing set. We investigate interpretability via coefficient interpretation for linear regression and lasso and variable importance for lasso and random forest.

The analysis was conducted using the statistical software R (R Core Team 2018). This document was written using the R packages `R Markdown` (Allaire et al. 2018), `knitr` (Xie 2018b), and `bookdown` (Xie 2018a). All materials used to conduct the analysis and compose the report can be found at <https://github.com/asbates/stat696/tree/master/reports/baseball>.

3 Analysis

3.1 Exploratory Analysis

There are a few areas of concern with the data use in this analysis. The first being a number of missing values for player salary². Almost 20% of the data has missing salary values. Most of the covariates were collected for the 1986 MLB season but salary was gathered at the beginning of the following season. Most of the missing salaries are likely due to retirement as 57 players retired in 1986³. The others might be missing a salary because the players returned to the minor leagues which is not uncommon in baseball. Regardless of the reason, with so many missing values we decided to impute them using k-nearest neighbors. The data set is not particularly large at 322 observations so omitting missing values would be leaving out a large chunk of the data.

Another potential issue with this data set is the correlations between the covariates. A plot of the correlation matrix for the continuous variables is given in the appendix (Figure 3). Several covariates have extremely high correlations, up to 98%. We also see groupings of covariates that have large correlations with each other. There are two groups of six variables each. One group

²None of the covariates had missing values.

³<http://www.baseball-almanac.com/yearly/final.php?l=NL&y=1986>

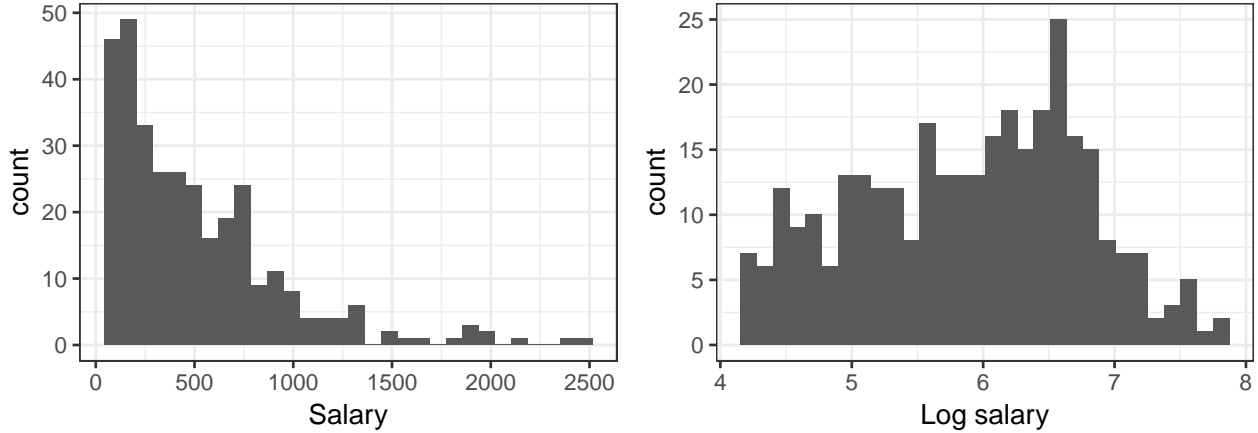


Figure 1: Histogram of player salary. The left plot is prior to transformation and the right plot is after log transforming.

containing career level variables and the other containing season level. For example, the number of hits and at bats have a correlation of 97%. This should not be surprising as players can only get a hit if they are at bat. But the more troubling issue is that both have high correlations with number of runs (92% for hits and 91% for at bats) and the number of runs batted in (81% for hits and 82% for at bats). Additionally, they both have moderate correlations with number of home runs and number of walks. These are likely to present problems, especially in the linear regression model. However, the model selection procedure is based on prediction error and fits subsets of the covariates so it may not include this group of variables. For now we refrain from removing any variables and reassess the issue after model selection.

Over half of the continuous variables exhibit skewness to some degree. This is of highest concern for the linear regression model but it may also affect the fit for lasso and random forest if most of the values are clumped together. Log transformations were used to account for skewness, favored for the interpretability of coefficients in the linear regression model. An example of this skewness is seen in Figure 1 where we have a histogram of salary (left) along with a histogram of log salary (right). The log transformation clearly helps with skewness but also note that after transformation salary is approximately normally distributed. Histograms of the remaining variables as well as any log transformed variables are provided in the appendix for reference. Also included are scatter plots of each covariate (or the log transformed covariate) versus log salary to assess the plausability of the linearity assumption for the linear regression model. The plots show that each predictor variable has an approximately linear relationship with log salary.

3.2 Modeling Fitting

After splitting the data into a testing and training set, each type of model was fit to the training data set with 5-fold cross-validation. The same cross-validation training and testing sets were used for each type of model. In each case, the final model was the one that resulted in the lowest root mean squared error (RMSE), averaged over the cross-validation runs. Prior to performing cross validation each continuous predictor was first centered by the mean and scaled by the standard deviation.

For the linear regression model a backwards stepwise procedure was used on each cross validation run. The resulting fit was the one with the lowest Akaike Information Criteria and this was then used to predict on the cross-validation hold out set. The model with the lowest mean cross-validation error was then evaluated. The coefficient for log put outs has a reasonably large p-value at 0.13 but the p-values for the remaining coefficients were satisfactory. The real concern with this model is the variance inflation factors (VIF). Number of at bats had a VIF of 19.2 and number of hits had a VIF of 16.9. These are quite large and indicate collinearity issues. This should not be surprising as we noted earlier that these variables are highly correlated (97%) and they both have moderate to high correlations with other variables. Number of at bats generally has higher correlations with other covariates than number of hits so it was decided to remove at bats and keep number of hits.

After removing number of at bats from the linear regression model the collinearity issue was mostly abated. The largest VIF for this model was moderate at 5.9 (number of runs batted in). Further diagnostics were done through plotting residuals versus fitted values and a QQ plot, both of which can be found in Figure 11 in the appendix. We see one concerning point with a rather large residual. The plot also shows a few points with variance a bit larger than most of the others. For the quantile-quantile plot, we see a point of concern that distances itself from the remaining points. However, recall that the primary purpose of this analysis is to compare predictive performance across each class of model and not necessarily to construct a model that perfectly satisfies its assumptions. Overall, both of the diagnostic plots in Figure 11 are reasonable approximations and are satisfactory for our purposes. Proceeding with this model, cross-validation was then performed using the same fit on each run in order to get an estimate of out-of-sample RMSE and allow for comparison with the other models.

Lasso and random forest were fit with ten values of their respective hyperparameters on each cross-validation iteration. The parameter for the lasso is the penalty parameter on the L_1 norm of the coefficients. Ten equally spaced values were used ranging from 0.01 to 1. The optimal parameter was 0.01. For the random forest the hyperparameter is the number of randomly selected predictor variables considered at each split. Values ranged from 2 to 19 with the optimal parameter found to be three.

3.3 Prediction Results

Table 1 displays the prediction error for each model across the cross-validation runs and on the testing set. As one would expect, the random forest model outperforms linear regression and lasso on both cross-validation error and test set error. The random forest is approximately 50% better than the alternatives in terms of relative error. The errors for linear regression and lasso are approximately the same which is a bit remarkable given that lasso often has improved predictive performance compared to linear regression (James et al. 2013, 203). Also of note, and perhaps most important, is the percent increase in error on the testing set compared to cross-validation error. Both linear regression and lasso saw an increase of about 30% while the random forest increase was only 7%.

Part of the performance gap between the linear methods and random forest is likely related to the groups of predictors with high correlations. Recall that there are two groups of six variables where each variable in a group has moderate to large correlation with the others. The optimal number of variables considered at each split for the random forest was three so it may be that only one or two predictors from each group is selected at a time. This would result in a better fit for each tree in

Table 1: Prediction results for linear regression, lasso, and random forest. Included are the mean cross-validation RMSE and test set RMSE. Also included is the increase in error from cross validation to testing and relative test set error.

Model	CV error	Test error	% Increase in error	Test relative error
Random Forest	0.43	0.46	7.3	1.00
Linear Regression	0.53	0.69	30.0	1.49
LASSO	0.54	0.71	32.0	1.53

the forest and hence a better fit when their predictions are averaged.

3.4 Model Interpretation

Interpretation is another aspect for which the three models considered here differ. Even if the main objective for building a model is prediction accuracy one usually wants to have an understanding of how the model is working. For example, we might be interested in how the covariates affect the response or which covariates are most important to the model. The type of model dictates what kind of inferences we can make. For linear regression, inferences most often take the form of interpreting the regression coefficients. The lasso is based on a standard linear regression model but the nature of the fitting process makes linear-regression-like inferences difficult. Random forests are similarly difficult to interpret as in linear regression. Although there has been some recent development on inference for both lasso (Lee et al. 2016) and random forest (Mentch and Hooker 2016), they beyond the scope of this analysis. Instead we take the traditional approach of model interpretation through feature importance measures.

A summary of the linear regression fit is given in Table 2. Included are coefficient estimates of the model along with their associated standard errors, p-values, and 95% confidence intervals. Because log transformations and centering and scaling of the covariates was performed, the interpretations of the coefficients are not as straightforward as when no transformations are performed. There are three types of predictors in the model and each warrants a slightly different interpretation. The simplest case is the covariates that were not log transformed, such as number of hits. Since the predictors were centered and scaled prior to fitting, we can interpret number of hits as follows. Hits has a standard deviation of 46 and its coefficient estimate is 0.18. So every 46 hits a player gets per season is associated with a $100 * 0.18$ or 18% increase in salary⁴. Now let's consider a covariate that was log transformed, log home runs. Log home runs has a standard deviation of 1 and a coefficient of 0.11. Suppose log home runs for a given player increases by one standard deviation. If we take $1 + 1$ and raise it to the power of 0.11 the result is 1.079. Subtract one from this and multiply by 100 and we get 7.9% which is the associated increase in salary we expect. The remaining type of predictor is the categorical variable for division. We leave this to the reader to interpret.

We should note that these inferences do not take into account the model selection procedure and as such should be handled with caution. Tibshirani et al. (2016) discuss how to account for this with forward stepwise regression but we use backwards selection here. In a higher stakes setting one should consider the implications of the model selection mechanism but for reasons of simplicity we assume our inferences are reasonable approximations.

⁴Of course, interpretations for linear regression must be made with the caveat that only the variable under question changes and the others are held fixed.

Table 2: Summary for linear regression model of log salary against hits, log home runs, runs batted in, walks, log put outs, log assists, log career runs batted in, and division. Coefficient estimates are provided along with their standard errors, p-values, and 95% confidence intervals.

Term	Estimate	SE	p-value	95% CI
Intercept	5.88	0.03	0.000	(5.82 , 5.94)
Hits	0.18	0.06	0.005	(0.06 , 0.31)
Log home runs	0.11	0.05	0.033	(0.01 , 0.21)
RBI's	-0.18	0.08	0.023	(-0.33 , -0.02)
Walks	0.08	0.04	0.052	(-0.00 , 0.17)
Log put outs	0.04	0.04	0.208	(-0.02 , 0.11)
Log assists	0.04	0.03	0.205	(-0.02 , 0.11)
Log career RBI's	0.59	0.04	0.000	(0.51 , 0.67)
Division: West	-0.09	0.03	0.005	(-0.16 , -0.03)

For lasso and random forest we interpret the models via variable importance measures. For variable importance of the lasso we use the absolute value of the regression coefficients. Feature importance for random forest is based on the difference in out-of-bag predictions before and after permuting a variable. Importance measures for both models are then scaled so that the most influential feature has an importance of 100. Plots of variable importance measures for lasso (left) and random forest (right) are given in Figure 2. Only variables with a positive importance are plotted. For the lasso the most important feature is log career walks. This may be because if a player can get on base, regardless of how they do so, they at least have a chance to score. The top 6 features for random forest are career level variables which seems reasonable because they take into account a players performance over their career instead of just a single season. Sometimes in baseball a players performance in a given season does not reflect their performance over their career; they may play better or worse this season compared to the past five. The random forest seems to be taking this into account. In contrast, the linear regression model only included one career level variable, log career runs batted in. Figure 2 also shows a varying difference across importance for lasso and random forest. For example, the difference between the top two features for the lasso is roughly 20 but for the random forest the difference is negligible. This may in part be due to the lasso only having 11 features compared to 18 in the random forest. Although we can not make the same type of statements about these models that we can in linear regression, variable importance measures still give us some insight into what is going on.

4 Conclusion

In this analysis we compared the ability of three classes of models to predict a baseball players salary given previous season information and overall career performance. The linear regression and lasso models achieved a comparable test set RMSE of 0.69 and 0.71, respectively. The random forest significantly outperformed the other models with a test set RMSE of 0.46. Random forest also less of a gap between its cross-validation and testing error compared to the other models. Linear regression and lasso had an approximately 30% increase in error while random forest had an increase of only 7%.

In terms of interpretability, the linear regression model has a slight advantage in that we can

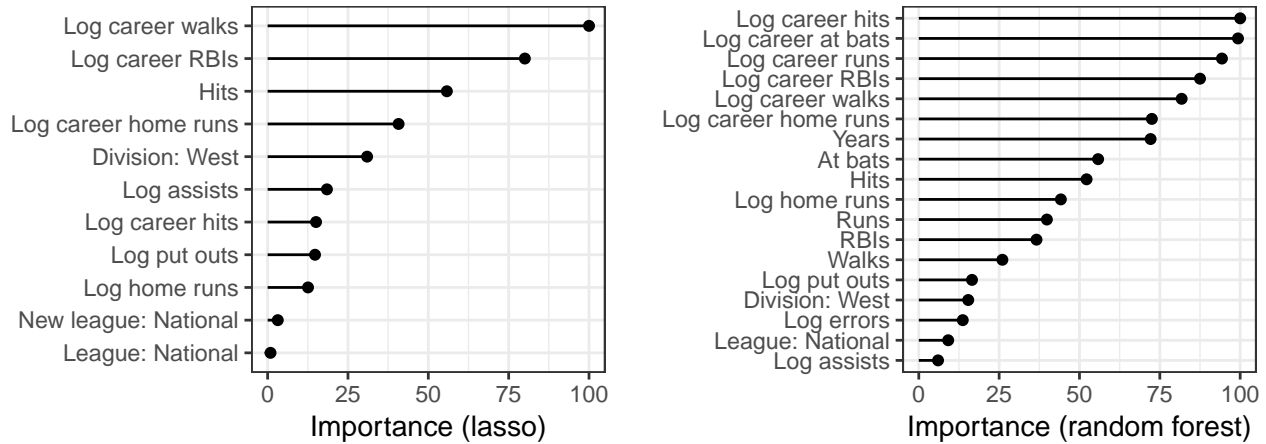


Figure 2: Variable importance plots for lasso (left) and random forest (right). Importance measures are scaled so that the most influential feature has an importance value of 100. Only variables with importance larger than 0 are included.

understand how a particular covariate relates to the response rather than just which covariates are important. That being said, we did not account for the linear regression model selection procedure so the inferences from the model are made with possible strong assumptions. Even though we did not make the same types of inferences for lasso and random forest, we nevertheless can obtain a high level understanding of the models through variable importance measures. There has been some recent results that allow for inferences in these models. They were beyond the scope of this analysis but perhaps a future study could investigate these methods.

It should be noted that this analysis does have its limitations. One issue is that the data is over 30 years old. Variables that have an impact on salary today are likely different than three decades ago. There have also been a substantial amount of new metrics since 1986 that may have more of an impact on salary than the simple ones included in this data⁵. It would be interesting to have a similar analysis performed on a contemporary data set to see if that is the case. The results might also have been different if a few more variables were included that would have been available in 1986. For example, player position could have allowed the models to treat pitchers different than outfielders. Alternative performance variables may have also mitigated the issues with correlations in this data set.

⁵See for example <http://m.mlb.com/glossary/advanced-stats>

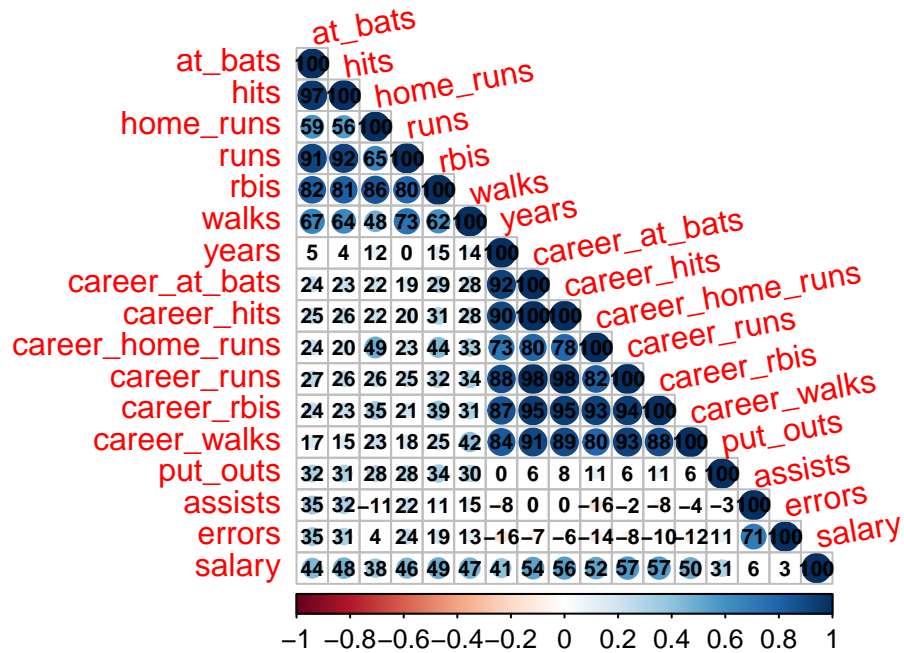


Figure 3: Correlation matrix plot of continuous variables prior to transformations. Correlations are displayed as a percent.

A Supplementary Figures

B R Code

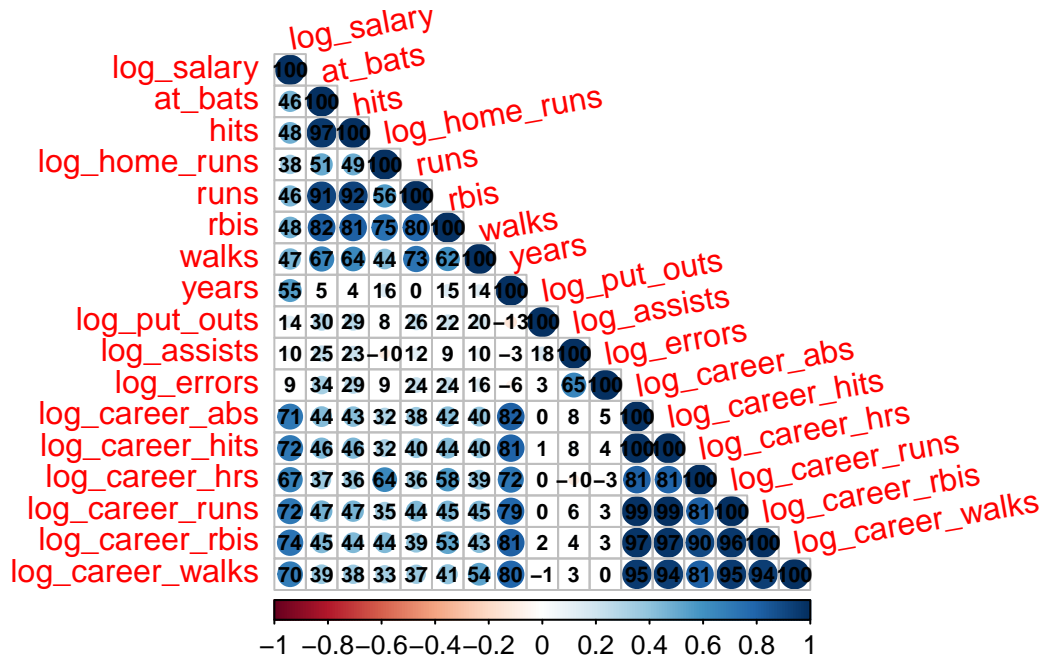


Figure 4: Correlation matrix plot of continuous variables after log transformations. Correlations are displayed as a percent.

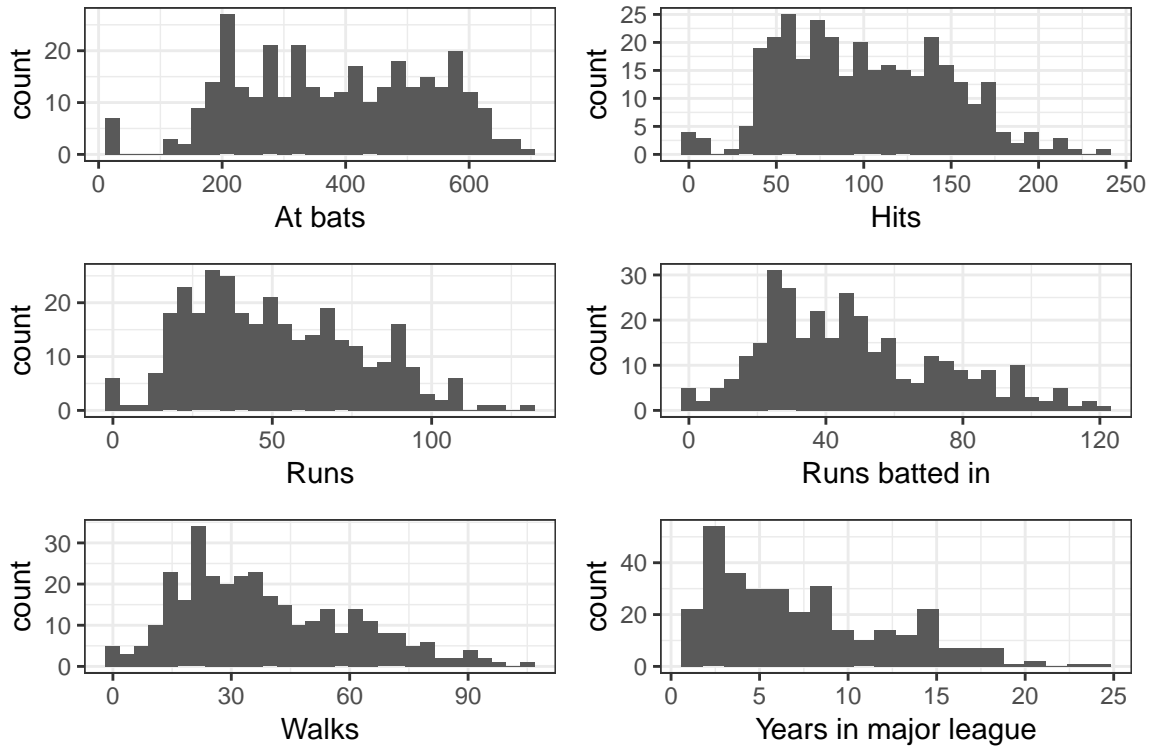


Figure 5: Histograms of covariates for which log transformations were not performed

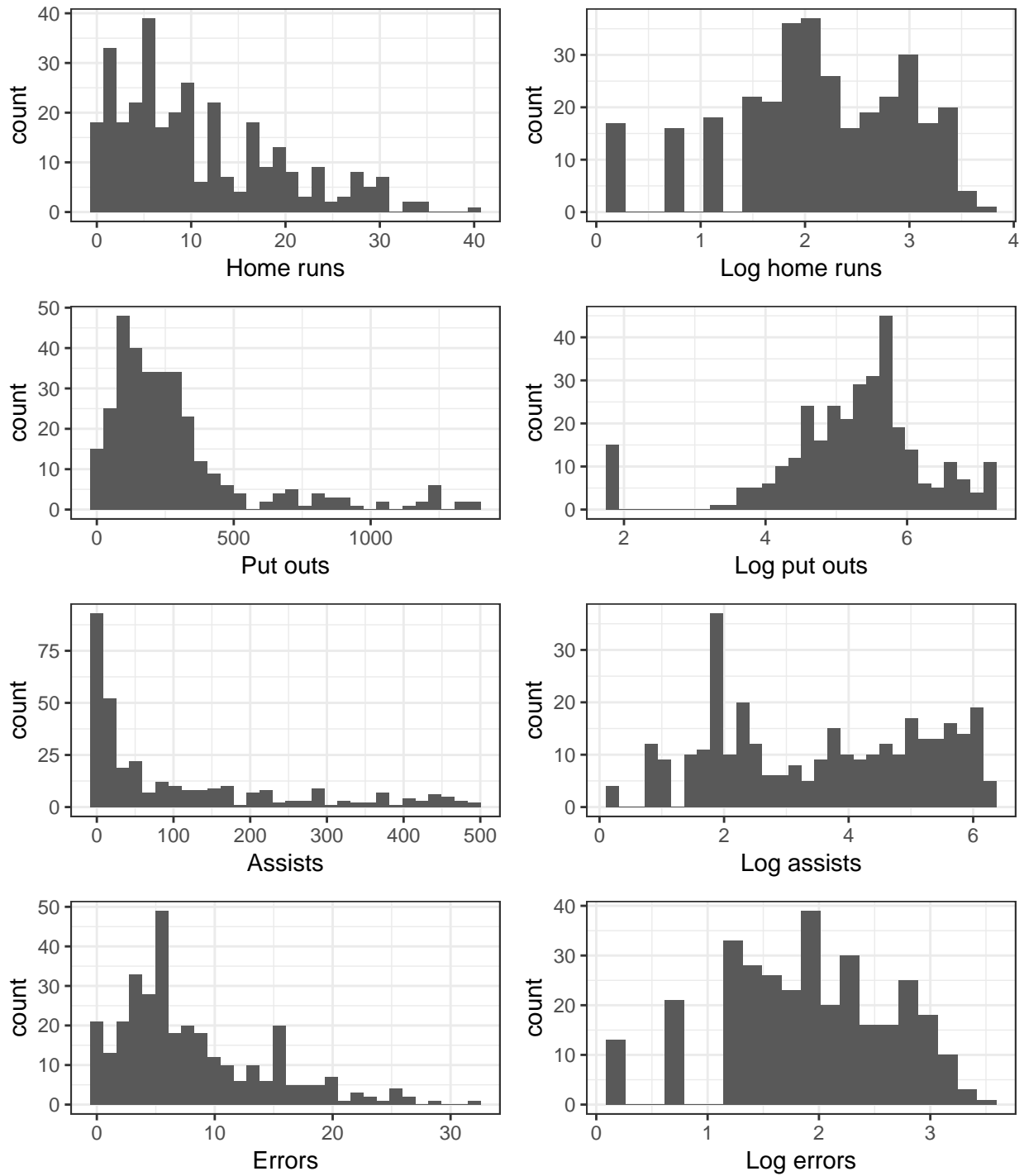


Figure 6: Histograms of season level covariates that were log transformed. The left side is the variable on the original scale and the right side is the variable on the log scale.

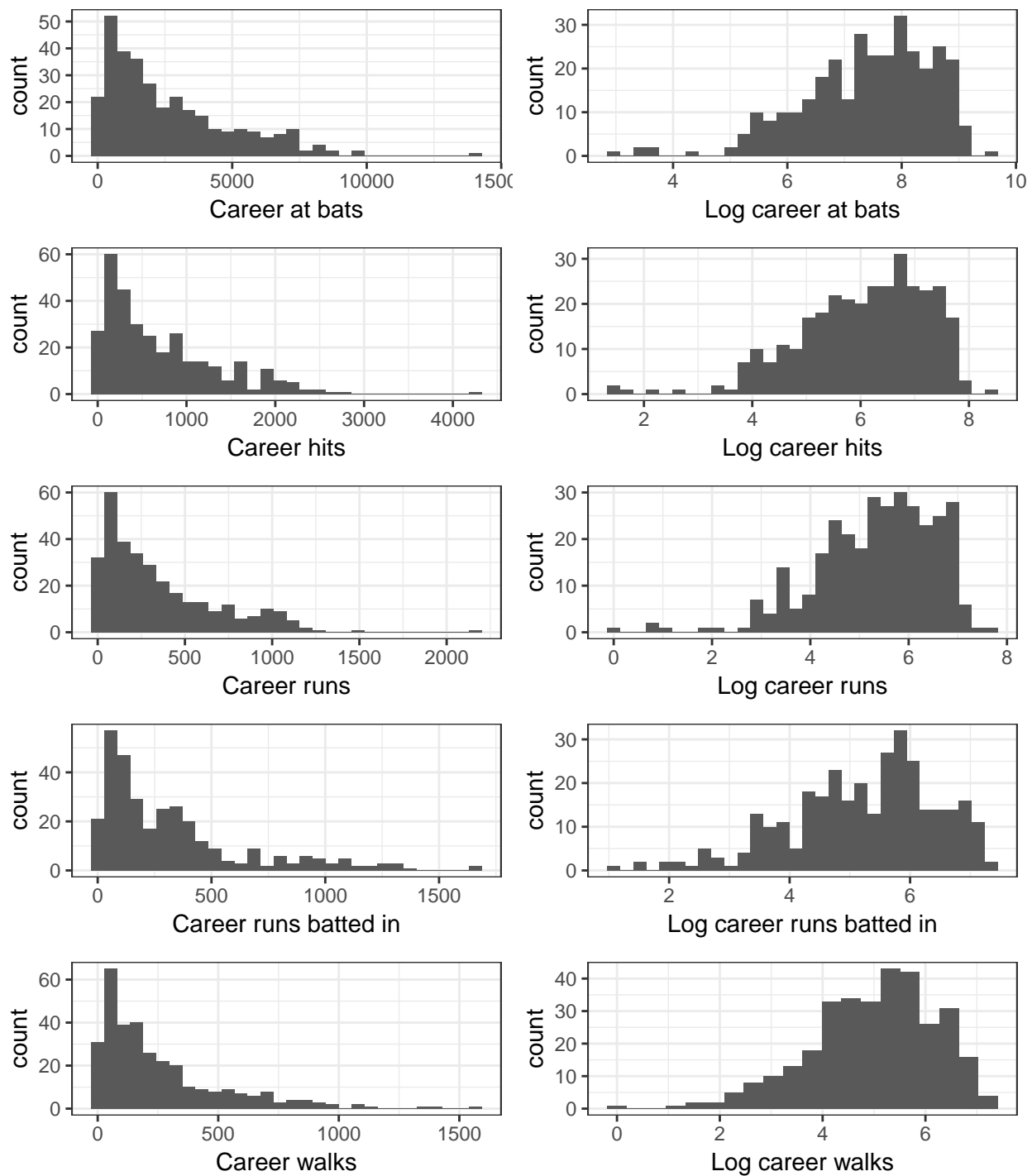


Figure 7: Histograms of career level covariates that were log transformed. The left side is the variable on the original scale and the right side is the variable on the log scale.

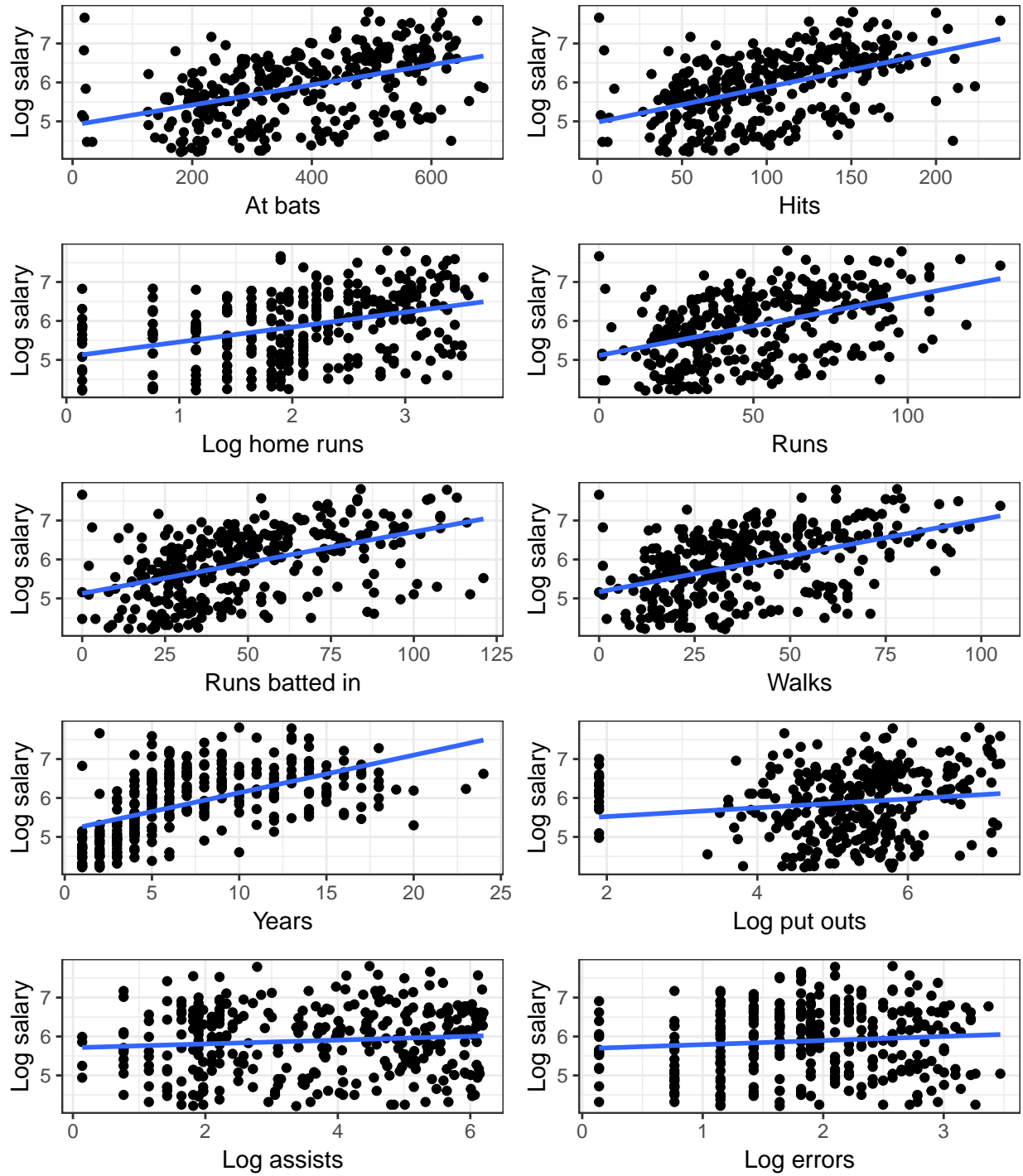


Figure 8: Scatter plots of log salary vs. season level covariates. The blue line is a smoothing line.

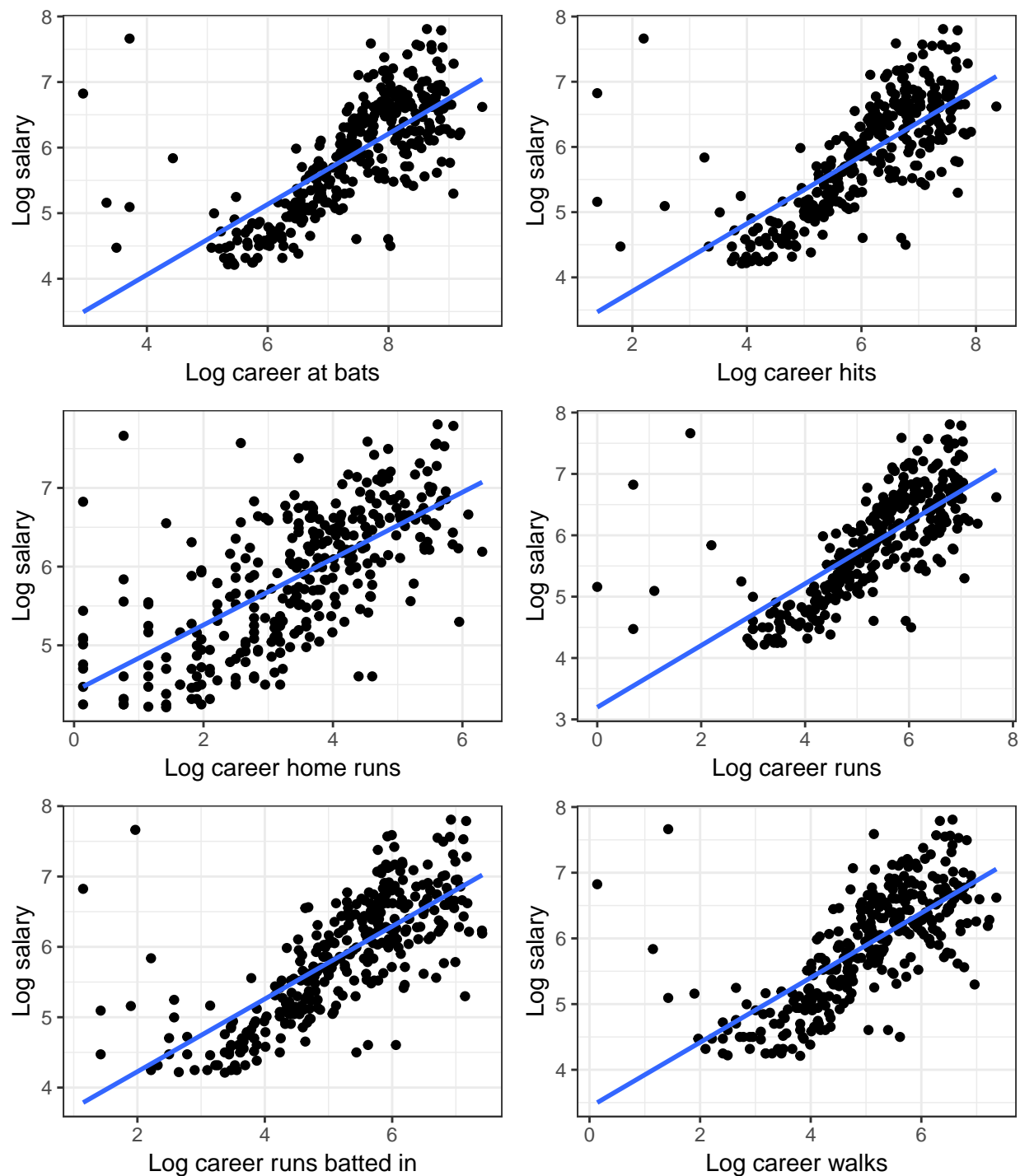
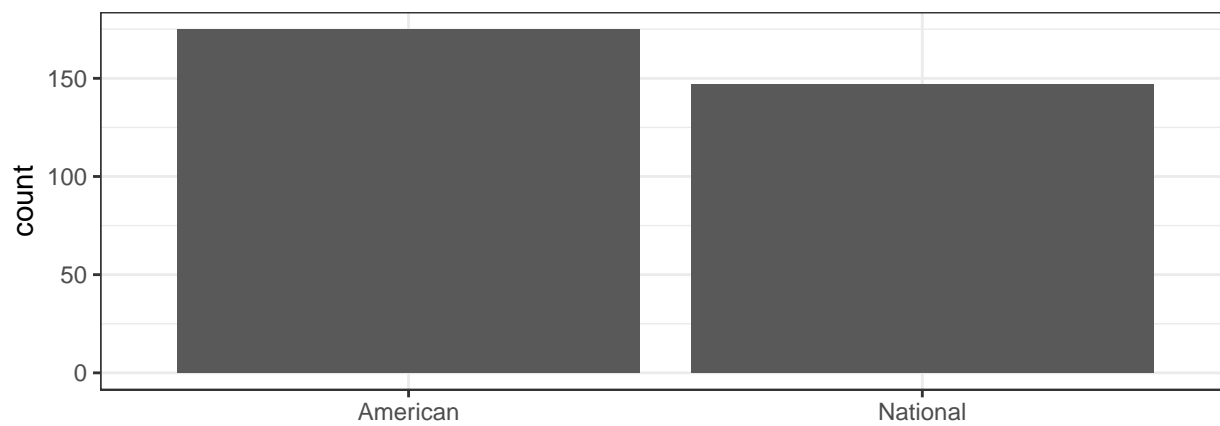
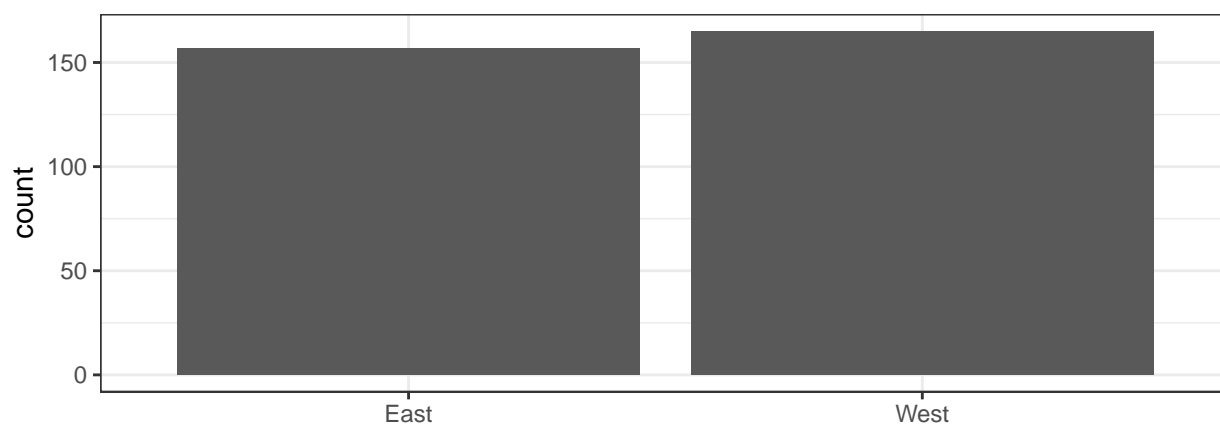


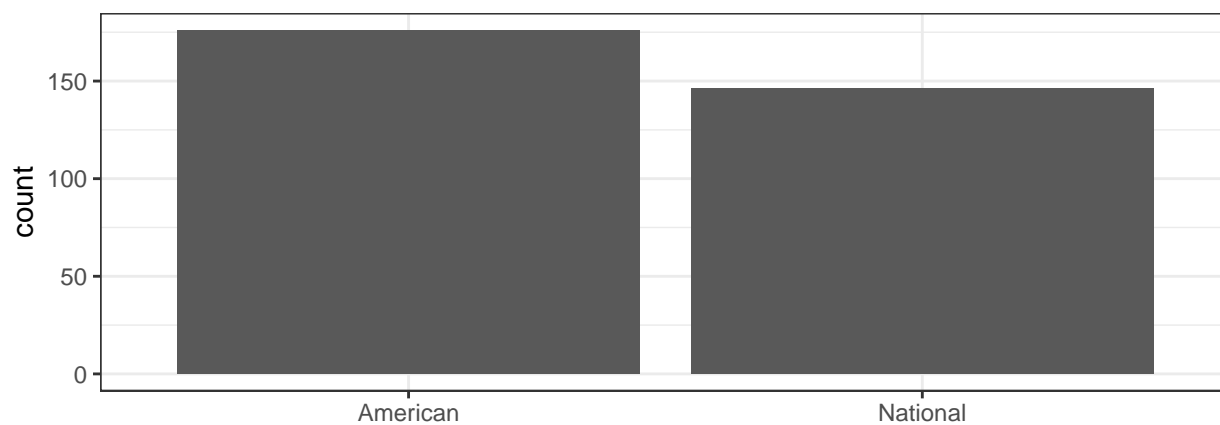
Figure 9: Scatter plots of log salary vs. career level covariates. The blue line is a smoothing line.



League at end of 1986 season



Division



League at start of 1987 season

Figure 10: Bar plots of categorical variables.

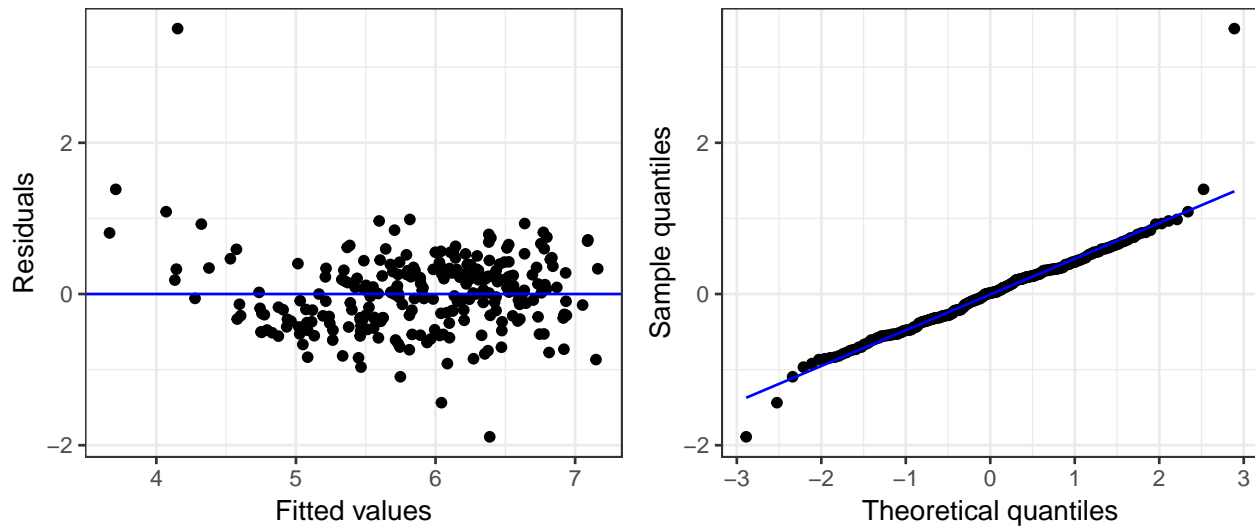


Figure 11: Diagnostic plots for the linear regression model. The left hand plot is residuals vs. fitted values. The right hand plot is a QQ plot of the residuals vs. normal quantiles. The blue lines are reference lines indicating zero (left) and the $y=x$ line (right).

References

- Allaire, JJ, Yihui Xie, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron Atkins, Hadley Wickham, Joe Cheng, and Winston Chang. 2018. *Rmarkdown: Dynamic Documents for R*. <https://CRAN.R-project.org/package=rmarkdown>.
- Breiman, Leo. 2001. “Statistical Modeling: The Two Cultures (with Comments and a Rejoinder by the Author).” *Statist. Sci.* 16 (3). The Institute of Mathematical Statistics: 199–231. doi:10.1214/ss/1009213726.
- Breiman, Leo, Adele Cutler, Andy Liaw, and Matthew Wiener. 2018. *RandomForest: Breiman and Cutler’s Random Forests for Classification and Regression*. <https://CRAN.R-project.org/package=randomForest>.
- Fox, John, Sanford Weisberg, and Brad Price. 2018. *Car: Companion to Applied Regression*. <https://CRAN.R-project.org/package=car>.
- Henry, Lionel, and Hadley Wickham. 2018. *Purrr: Functional Programming Tools*. <https://CRAN.R-project.org/package=purrr>.
- James, Gareth, Daniela Witten, Trevor Hastie, and Rob Tibshirani. 2013. *An Introduction to Statistical Learning With Applications in R*. Springer.
- . 2017. *ISLR: Data for an Introduction to Statistical Learning with Applications in R*. <https://CRAN.R-project.org/package=ISLR>.
- Jed Wing, Max Kuhn. Contributions from, Steve Weston, Andre Williams, Chris Keefer, Allan Engelhardt, Tony Cooper, Zachary Mayer, et al. 2018. *Caret: Classification and Regression Training*.

<https://CRAN.R-project.org/package=caret>.

Kuhn, Max, and Hadley Wickham. 2018. *Recipes: Preprocessing Tools to Create Design Matrices*. <https://CRAN.R-project.org/package=recipes>.

Lee, Jason D., Dennis L. Sun, Yuekai Sun, and Jonathan E. Taylor. 2016. “Exact Post-Selection Inference, with Application to the Lasso.” *Ann. Statist.* 44 (3). The Institute of Mathematical Statistics: 907–27. doi:10.1214/15-AOS1371.

Mentch, Lucas, and Giles Hooker. 2016. “Quantifying Uncertainty in Random Forests via Confidence Intervals and Hypothesis Tests.” *Journal of Machine Learning Research* 17 (26): 1–41. <http://jmlr.org/papers/v17/14-168.html>.

Müller, Kirill. 2017. *Here: A Simpler Way to Find Your Files*. <https://CRAN.R-project.org/package=here>.

R Core Team. 2018. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.

Ripley, Brian. 2018. *MASS: Support Functions and Datasets for Venables and Ripley’s Mass*. <https://CRAN.R-project.org/package=MASS>.

Tibshirani, Ryan J., Jonathan Taylor, Richard Lockhart, and Robert Tibshirani. 2016. “Exact Post-Selection Inference for Sequential Regression Procedures.” *Journal of the American Statistical Association* 111 (514). Taylor & Francis: 600–620. doi:10.1080/01621459.2015.1108848.

Wei, Taiyun, and Viliam Simko. 2017. *Corrplot: Visualization of a Correlation Matrix*. <https://CRAN.R-project.org/package=corrplot>.

Wickham, Hadley, Winston Chang, Lionel Henry, Thomas Lin Pedersen, Kohske Takahashi, Claus Wilke, and Kara Woo. 2018. *Ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics*. <https://CRAN.R-project.org/package=ggplot2>.

Wickham, Hadley, Romain François, Lionel Henry, and Kirill Müller. 2018. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.

Xie, Yihui. 2018a. *Bookdown: Authoring Books and Technical Documents with R Markdown*. <https://CRAN.R-project.org/package=bookdown>.

———. 2018b. *Knitr: A General-Purpose Package for Dynamic Report Generation in R*. <https://CRAN.R-project.org/package=knitr>.

Zhu, Hao. 2018. *KableExtra: Construct Complex Table with ‘Kable’ and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.