# **Analyzing the Relationship Between Bike Shares** and Seasonal Conditions

Data Analysis Project Report - STAT 420, Summer 2021, D. Unger

- Introduction
- Methods • Results
- Discussion
- Appendix

# Introduction

In this report, we intend to model bike share frequency in London based on seasonal conditions (e.g. weather conditions, whether it's a holiday, etc.) for purposes of prediction. We are creating this model because these insights would be useful to bike/scooter sharing companies and transport authorities. We also find this data ideal for exploring & applying the techniques we've learned in STAT 420 whilst remaining easily interpretable.

### **Report Layout**

In this report, the Methods section will offer a narrative of the step-by-step decision-making process we followed to arrive at our final model. In the Results section, we will show numerical and graphical summaries of the final model. In the Discussion section, we will discuss our results and frame them in the context of the data.

#### **Dataset Background Information**

The data set we chose is the London bike sharing dataset from Kaggle. It consists of the number of new rentals (shares) of bikes owned by the London transport authority (Transport for London). Counts are provided for (nearly) every hour between January 4th, 2015 and January 3rd, 2017. This count is accompanied by weather & date information pertinent to the time the count was recorded (e.g. humidity, wind speed, whether it's a holiday, current season, etc.)

The data set is a data frame with 10 variables and 17,414 observations.

Variable **Description** 

timestamp	Timestamp (hourly)
cnt	Number of new bike shares/rentals
t1	Real temperature (°C)
t2	Relative temperature (°C)
hum	Relative humidity (%)
wind_speed	Wind speed (km/h)
weather_code	Weather type $(1 = \text{Clear}, 2 = \text{Partly cloudy}, 3 = \text{Mostly cloudy}, 4 = \text{Overcast}, 7 = \text{Rain}, 10 = \text{Thunderstorm}, 26 = \text{Snow})$
is_holiday	Is a holiday ( $0 = \text{non-holiday}, 1 = \text{holiday}$ )
is_weekend	Is the weekend (O = weekday, 1 = weekend)
season	Meteorological season (0 = spring, 1 = summer, 2 = fall, 3 = winter)

# **Methods**

In this section, we will provide a narrative step-by-step decision-making process throughout our analysis as we adjust the model and attempt to validate model assumptions. Below, we will make use of strategies such as multiple linear regression, dummy variables, interaction terms, residual diagnostics, outlier diagnostics, transformations of the response, and model selection.

#### **Read and Clean the Dataset**

A description of the dataset can be found in the Introduction. We converted categorical variables to factors, dropped the timestamp (unusable), and dropped an observation with cnt = 0 to prevent issues with Box-Cox later.

```
ds = read.csv("london merged.csv")
ds$weather code = as.factor(ds$weather code)
ds$is holiday = as.factor(ds$is holiday)
ds$is weekend = as.factor(ds$is weekend)
ds$season = as.factor(ds$season)
ds = ds[, setdiff(names(ds), c("timestamp"))] # drop `timestamp` column
ds = ds[ds$cnt > 0, ] # drop 1 observation with `cnt` = 0
```

#### **Perform Train-test Split**

We split the data into a train set & test set so we can get an estimate of model performance on unseen data.

set.seed(19990420) ds idx = sample(1:nrow(ds), (0.80)\*nrow(ds)) ds trn = ds[ds idx,] ds\_tst = ds[-ds\_idx,]

#### **Check Pairwise Correlations**

Here, we check the pairs plot and correlation matrix for highly correlated variables. We remove t1 because it's highly correlated with t2.

pairs(ds trn, col="dodgerblue")



ds trn = ds trn[setdiff(names(ds trn), c("t1"))] # drop `t1` column

#### Two-Way Interaction, Additive Model, and Backwards AIC Search

We begin by fitting an additive and two-way interaction model. We try a backwards AIC search on both. The model found via backwards AIC search on the additive model has a higher adjuted  $R^2$  so we move forward with it

backwards AIC search on the additive model has a higher adsuted A <sup>-</sup> , so we move forward with it.
<pre>fit_additive = lm(cnt ~ ., data = ds_trn) fit_add_back_aic = step(fit_additive, direction = "backward", trace = 0) summary(fit_add_back_aic)\$adj.r.squared</pre>
## [1] 0.2888
fit interaction = $lm(cnt \sim . ^ 2, data = ds trn)$
fit int back aic = step(fit interaction, direction = "backward", trace = $0$ )
summers (fit int heat sic) (add a summer d
Summary(IIt_Int_Dack_aic)\$adj.r.squared
## [1] 0.3353

#### **Drop Insignificant Variables From Interaction Model**

Here, we analyze the significance of the predictors. We decide to drop windspeed and is holiday because they have the lowest Pr(>|t|). We fit a new model with this reduced set of predictors.

summary(fit\_int\_back\_aic)\$call

<pre>## lm(formula = cnt ~ t2 + hum + wind_speed + weather_code + is_holiday + ## is_weekend + season + t2:hum + t2:is_holiday + t2:is_weekend + ## t2:season + hum:wind_speed + hum:weather_code + hum:is_holiday + ## hum:is_weekend + hum:season + wind_speed:weather_code + wind_speed:is_holiday + ## wind_speed:is_weekend + wind_speed:season + weather_code:is_holiday + ## weather_code:is_weekend + weather_code:season + is_holiday:season + ## is_weekend:season, data = ds_trn)</pre>						
<pre>summary(fit_int_base)</pre>	ack_aic)\$coef	f[1:15, ] # windspeed and is_holiday are insignificant				
##	Estimate {	Std. Error t value Pr(> t )				
## (Intercept)	1727.5835	188.441 9.1678 5.500e-20				
## t2	88.1914	9.431 9.3512 9.973e-21				
## hum	-17.9135	2.340 -7.6554 2.054e-14				
## wind speed	0.6677	5.934 0.1125 9.104e-01				
## weather code2	54.5473	145.761 0.3742 7.082e-01				
## weather code3	-86.4679	172.991 -0.4998 6.172e-01				
## weather code4	-2210.6179	263.564 -8.3874 5.446e-17				
## weather code7	351.9950	245.921 1.4313 1.524e-01				
## weather code10	1126.0745	4056.305 0.2776 7.813e-01				
## weather_code26	415.6930	3373.443 0.1232 9.019e-01				
## is_holiday1	949.2060	580.054 1.6364 1.018e-01				
## is_weekend1	946.1718	139.060 6.8040 1.059e-11				
## season1	-377.8921	220.697 -1.7123 8.687e-02				
## season2	781.8681	198.463 3.9396 8.201e-05				
## season3	264.9122	199.248 1.3296 1.837e-01				

fit rm insig = lm(cnt ~ t2 + hum + weather code + is weekend + season + t2:hum + t2:is weekend + t2:season + hum:weather code + hum:is weekend + hum:season + weather code:is weekend + weather c ode:season + is weekend:season, data = ds trn)

#### **Checking Model Assumptions & Performance**

Here, we check for model assumptions & overfitting. We find that the LOOCV-RMSE is infinite (likely due to high-leverage points). We also find the adjusted  $R^2$  for comparison. The Fitted vs Residuals shows that the constant variance and linearity assumptions are suspect. The Q-Q Plot shows that the normality assumption is suspect.



**Fitted versus Residuals** 



### **Box-Cox Transformation of the Response**

Here, we try a box-cox transformation of the response variable.

bc = boxcox(fit rm insig, plotit = TRUE, lambda = seq(0.2, 0.35, by = 0.01))

(lambda = bc\$x[which.max(bc\$y)])

## [1] 0.2712

fit boxcox int = lm(((cnt ^ lambda - 1) / lambda) ~ t2 + hum + weather\_code + is\_weekend + seaso n + t2:hum + t2:is weekend + t2:season + hum:weather code + hum:is\_weekend + hum:season + weathe r\_code:is\_weekend + weather\_code:season + is\_weekend:season, data = ds\_trn)

4

#### **Checking Model Assumptions & Performance**

We see that the interaction model with the Box-Cox transformation has several (32) variables with a high VIF and infinite LOOCV RMSE. To lower the VIF of the variables and get a defined LOOCV RMSE, we will next try a model without any interaction terms. Too, our normality assumption looks substantially improved & the Fitted Versus Residuals plot shows slight improvement towards the linearity assumption.

c(sum(vif(fit boxcox int) > 5), max(vif(fit boxcox int)))

## [1] 32 2617

calc loocv rmse(fit boxcox int)

## [1] Inf

summary(fit boxcox int)\$adj.r.squared

## [1] 0.3742

20

10

0

-10

## [1] 0

5

Residuals

plot diagnostics(fit boxcox int)

**Fitted versus Residuals** 



#### **Remove Interaction Terms and Compare**

sum(vif(fit\_boxcox\_add) > 5)

10

15

Here, we modify the Box-Cox model to remove interaction terms. We prefer a simpler model where possible. This new model has a good LOOCV RMSE and none of its predictors have a high VIF (> 5).

<pre>calc_loocv_rmse(fit_boxcox_add)</pre>
## [1] 5.688
Checking Model Assumptions & Performance
Here, we see the same possible violation of the linearity & constant variance assumptions with a decent normality assumption adherence.
<pre>summary(fit_boxcox_add)\$adj.r.squared</pre>
## [1] 0.344
<pre>plot_diagnostics(fit_boxcox_add)</pre>
High leverage, Outlier, and Influential Observations
Here, we check for high leverage, outlier, and influential observations. We see that about 4% of the observations are influential (high leverage AND large residual).
unusual_observations(fit_boxcox_add)
<pre>## [1] "High Leverage: 62 points, 0.45% of points" ## [1] "Outliers: 654 points, 4.69% of points" ## [1] "Influential: 435 points, 3.12% of points"</pre>

### **Remove influential points**

Here, we refit the Box-Cox additive model after removing influential observations. The coefficients change but not drastically.

cooks dists = cooks.distance(fit boxcox add) fit boxcox rm infl = lm(formula = ((cnt ^ lambda - 1) / lambda) ~ hum + season + weather code + t2 + is weekend, data = ds trn, subset = cooks dists <= 4 / length(cooks dists)) coef(fit boxcox rm infl)

##	(Intercept)	hum	season1	season2	season3
##	29.9074	-0.2292	-1.0325	1.2067	1.6361
##	weather_code2	weather_code3	weather_code4	weather_code7	weather_code10
##	2.4961	2.8226	-0.7725	0.7407	-0.4841
##	weather_code26	t2	is_weekend1		
##	-1.5231	0.3160	-0.4829		

#### **Check Model Assumptions and Performance**

Here, we see a greatly improved Adjusted  $R^2$ , still good VIFs, and a good LOOCV RMSE. We see that this model has the best adjusted  $R^2$  yet. The model still doesn't quite follow the constant variance or linearity assumption, but the normality assumption does seem to be satisfied.





#### High Leverage, Outlier, and Influential Observations

Here, we see the proportion of influential points has improved greatly and the percent of outliers and high leverage points has also been reduced.

unusual_observations(fit_boxcox_rm_infl)					
## [1]	"High Leverage:	27	points,	0.2% of points"	
## [1]	"Outliers:	603	points,	4.47% of points"	
## [1]	"Influential:	282	points,	2.09% of points"	

#### **Compare the Models**

Finally, we compare all of the models we've fit using the adjusted  $R^2$ , number of coefficients, and Test RMSE (using heldback testing data). We see that the final model (Box-Cox interaction with removed influential points), is a good tradeoff. See the next sections for further discussion.

compare\_models(models = list(fit\_additive = fit\_additive, fit\_add\_back\_aic = fit\_add\_back\_aic, f it interaction = fit interaction, fit int back aic = fit int back aic, fit boxcox int = fit boxc ox\_int, fit\_boxcox\_add = fit\_boxcox\_add, fit\_boxcox\_rm\_infl = fit\_boxcox\_rm\_infl), lambda = lam bda)

Model Name	Test RMSE	Adj. R^2	# Coefs
fit_additive	914.3	0.2888	15
fit_add_back_aic	914.3	0.2888	15
fit_interaction	884.4	0.3350	88
fit_int_back_aic	884.2	0.3353	80
fit_boxcox_int	928.3	0.3742	55
fit_boxcox_add	950.2	0.3440	13
fit_boxcox_rm_infl	958.8	0.4186	13

### **Results**

The previous section (Methods) discussed our approach to finding a good model for our dataset. Readers should refer to it for the step-by-step narrative and intermediate models. The model we settled on was of the form (cnt ^ 0.2712 - 1) / 0.2712 ~ hum + season + weather code + t2 + is weekend.

**Final Model Model Assumptions and Performance** 

Our final model has the highest adjusted  $R^2$  we could find. It also has good VIFs and a good LOOCV-RMSE. The final model, like the others doesn't seem to adhere to constant variance or linearity assumptions, but the normality assumption does seem to be satisfied.

## [1] 0.41	86				
##	hum	season1	season2	season3 weather code2	



## [1] 5.224

**Fitted versus Residuals** 



Normal Q-Q Plot

**Final Model High Leverage, Outlier, and Influential Observations** 

Our final model has low proportions of high leverage, outlier, and influential observations.

## [	[1]	"High Leverage:	27	points,	0.2% of points"
## [	[1]	"Outliers:	603	points,	4.47% of points"
## [	[1]	"Influential:	282	points,	2.09% of points"

### **Final Model Comparison**

To further justify our choice of final model, we compare all of the models we've fit using the adjusted  $R^2$ , number of coefficients, and Test RMSE (using held-back testing data). We see that the final model (Box-Cox interaction with removed influential points) has the best adjusted  $R^2$  while maintaining a good Test RMSE and a very low number of coefficients. Too, like the other Box-Cox transformation models, our final model aligns better with the normality assumption. For all of these reasons, we believe the fit boxcox rm infl model is the best of the many models we investigated.

Model Name	Test RMSE	Adj. R^2	# Coefs
fit_additive	914.3	0.2888	15
fit_add_back_aic	914.3	0.2888	15
fit_interaction	884.4	0.3350	88
fit_int_back_aic	884.2	0.3353	80
fit_boxcox_int	928.3	0.3742	55
fit_boxcox_add	950.2	0.3440	13
fit_boxcox_rm_infl	958.8	0.4186	13

# Discussion

In the Methods section, we offered a narrative of the step-by-step decision-making process we followed to arrive at our final model. In the results section, we showed numerical and graphical summaries of the final model. In this section, we will discuss our results and frame them in the context of the data.

The final model we settled on was of the form

 $(cnt ^ 0.2712 - 1) / 0.2712 \sim hum + season + weather_code + t2 + is_weekend. It has the highest adjusted <math>R^2$  we could find. It also has good VIFs and a good LOOCV-RMSE. The final model, like the others doesn't seem to adhere to constant variance or linearity assumptions, but the normality assumption does seem to be satisfied. An addition, we compared all of the models we've fit using the adjusted  $R^2$ , number of coefficients, and Test RMSE (using held-back testing data). We see that the final model (Box-Cox interaction with removed influential points) has the best adjusted  $R^2$  while maintaining a good Test RMSE and a very low number of coefficients. Too, like the other Box-Cox transformation models, our final model aligns better with the normality assumption. For all of these reasons, we believe the fit boxcox rm infl model is the best of the many models we investigated.

The original goal of this model is to predict bike share frequency in London based on seasonal conditions (e.g. weather conditions, whether it's a holiday, etc.) for purposes of prediction. Our final model could offer useful insights to bike/scooter sharing companies and transport authorities. This model also helped us explore & apply the techniques learned in STAT 420 whilst remaining easily interpretable.

We believe our model achieves this goal as it was able to perform with a 959 RMSE on a held-out test dataset. It also comes closest to meeting LINE assumptions. Because it has few coefficients and no interaction terms, it remains easy to interpret and understand. Too, it has the highest adjusted  $R^2$  of the models we tried.

So, we are satisfied with the final model we've found and believe it would prove useful in the real-world problem of predicting bike share frequency based on seasonal conditions.

# Appendix

## **Helper Functions**

**Diagnostics Plots:** 

This function plots the Fitted versus Residual and Q-Q Plot for a model.

```
plot diagnostics = function(model, pcol = "grey", lcol = "dodgerblue") {
  Layout(matrix(c(1,2), nrow = 1, ncol = 2, byrow = TRUE))
  par(mfrow=c(1,2))
 plot(fitted(model), resid(model), col = pcol, xlab = "Fitted", ylab = "Residuals", main = "Fit
ted versus Residuals")
  abline(h = 0, col = lcol)
  qqnorm(resid(model), col = pcol, main = "Normal Q-Q Plot")
  qqline(resid(model), col = lcol)
}
```

#### **Unusual Observations**

This function prints out the number and proportion of high leverage, outlier, and influcential observations.

```
unusual observations = function(model) {
   leverages = hatvalues(model)
   print(paste("High Leverage: ", sum(leverages > 2*mean(leverages)), " points, ", round(mean(lev
 erages > 2*mean(leverages)) * 100, 2), "% of points", sep = ""))
   rstandard_abs = abs(rstandard(model))
   print(paste("Outliers: ", sum(rstandard_abs > 2), " points, ", round(mean(rstandard_abs > 2))
 2) * 100, 2), "% of points", sep = ""))
   cooks_dists = cooks.distance(model)
   cd_cutoff = 4 / length(cooks_dists)
   print(paste("Influential: ", sum(cooks_dists > cd_cutoff), " points, ", round(mean(cooks_dis
 ts > cd_cutoff) * 100, 2), "% of points", sep = ""))
 }
Compare Models
This function outputs a table comparing models using the test RMSE, adjusted R^2, and number of betas.
 compare models = function(models, lambda) {
   results = data.frame(name = character(), test_rmse = numeric(), adj r_squared = numeric(), num
 betas = numeric())
   for (name in names(models)) {
     model = models[[name]]
     pred = predict(model, newdata = ds_tst)
     test_rmse = ifelse(grepl("boxcox", name, fixed = TRUE), rmse(ds_tst$cnt, (pred * lambda + 1
 ) ^ (1 / lambda)), rmse(ds_tst$cnt, pred))
     adj r squared = summary(model)$adj.r.squared
     num betas = length(coef(model))
     results[nrow(results) + 1, ] = list(name, test rmse, adj r squared, num betas)
   }
```

kable(results, col.names = c("Model Name", "Test RMSE", "Adj. R^2", "# Coefs"))

**Group Members** 

}

- Darci Peoples (darciap2)
- Alex Koczwara (alexk3) Yash Patel (ypatel42)