

# **Predicting the Present with Google Trends**

## Forecasting Flu Season

STA 237A Final Project

Kenneth Broadhead, Xuezhen Li  
December 13, 2019

# 1 Abstract

In this paper, we consider using google trends search data to predict current, but presently unknown values, of time series data of interest. Two examples are discussed, an economic example involving automobile sales, and a public health example involving seasonal flu data collected from the CDC. Potential benefits and limitations of this method of “forecasting the present” are also discussed.

# 2 Introduction

Economic reports for a given time period are plagued by a significant reporting lag. Similarly, weekly flu Data from the CDC’s network of Hospitals and Clinics takes time to compile and release to the public. Understandably, economic researchers and public health officials, respectively, would like to have access to these reports much sooner. However, the reporting delay in both cases is due to the fact that these data sets are large and complex, taking time compile and publish. Reporting delay is often simply unavoidable.

However, Google Trends data for search queries related to the information contained in these reports is available in almost real time. It is not unreasonable to suspect that these search queries are correlated with the subject matter of the reports of interest, for the search data is reflective of the actions, (or experiences) of those about whom the reports are in some sense about. We therefore endeavor to investigate whether it is possible to predict the current, but yet to be released, values of these reports using current Google Trends data. We first explain to the reader relevant aspects of Google Trends, we then explore two examples: predicting automobile sales, and predicting the national prevalence of flu in the US.. We close with a discussion of the potential benefits and limitations of this method of ”predicting the present,” as well as possible methods for improvement that should be investigated further.

# 3 Google Trends

Google Trends provides a time series index of the volume of queries users enter into Google in a given geographic area. Figure 1 shows the search index for [UC Davis] in the United States in the past five years (As a side note, we the sudden spike in interest in UC Davis occurred during March madness, when UC Davis competed in the NCAA men’s basketball tournament). The maximum query share in the time period specified is normalized to be 100 and the query share at the initial date being examined is normalized to be zero. In addition to the region and period of time, we can also restrict the search index to a specific category. Google classifies search queries into about 30 categories at the top level and about 250 categories at the second level using a natural language classification engine.

The last thing we mention, is the limitation of acquiring Google Trends data. The longer the time period you choose, the longer the time interval between two consecutive search indexes. For example, if the time period is longer than five years, you can only get the monthly Google trends data instead of weekly data.

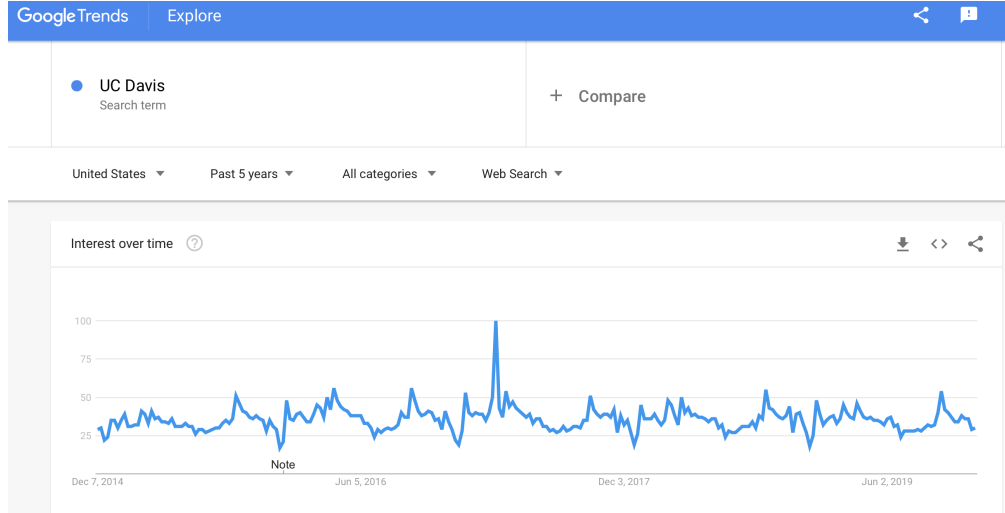


Figure 1: Search index for [UC Davis] in the United States

## 4 Economic Data Example

As the first economic data example we use the “Motor Vehicles and Parts Dealers” series from the U.S. Census Bureau “Advance Monthly Sales for Retail and Food Services” report.[1] Figure 2 shows the data from 2005 to 2011 after the log transformation.

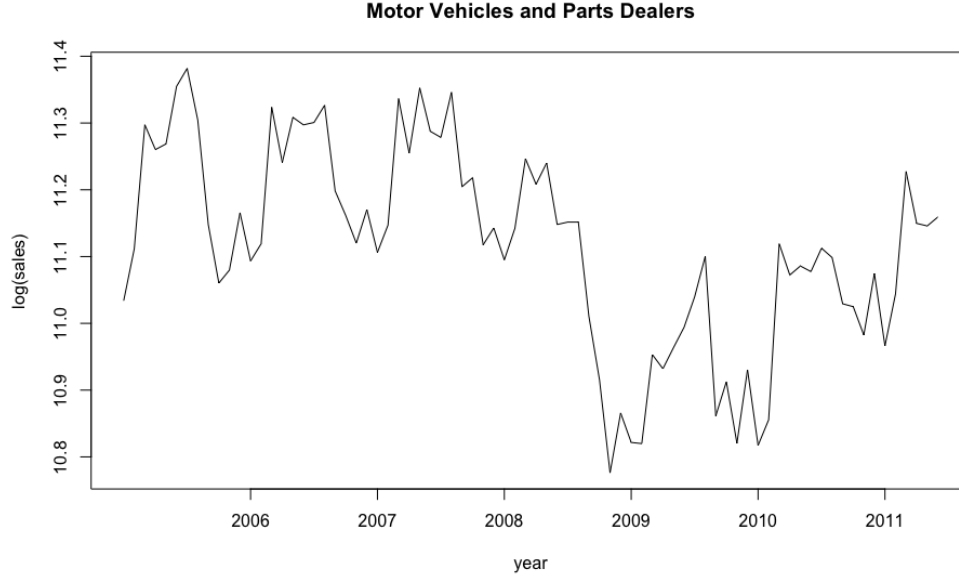


Figure 2:  $\log(\text{sales})$  of Motor Vehicles and Parts

By looking at the ACF plot we decide to choose AR(5) seasonal model as the initial full model and use all subset regression method to implement model selection. Figure 3 shows that the AR(1) seasonal model is the best model based on BIC and Mallows'  $C_p$  criteria. The  $R^2$  of the AR(1) seasonal model is 0.7211.

lag1	lag2	lag3	lag4	lag5	lag12	R <sup>2</sup>	adj R <sup>2</sup>	Cp	BIC
*						0.6626	0.6581	15.9899	-76.0299
*					*	0.7211	0.7137	2.3798	-86.5332
*			*		*	0.7246	0.7135	3.4406	-83.1677
*	*		*		*	0.7328	0.7182	3.2638	-81.1579
*	*	*	*		*	0.7336	0.7151	5.0420	-77.0443
*	*	*	*	*	*	0.7338	0.7113	7.0000	-72.7337

Figure 3: all subset regression

The base model is

$$y_t = b_0 + b_1 y_{t-1} + b_{12} y_{t-12} + e_t \quad (\text{base model})$$

Google Trends contains several automotive-related categories. We add two automotive-related categories from the Google trends.[1] One is Trucks and SUVs, the other one is Auto Insurance. To check the temporal aspects of the relationship between response and google trends covariates, we plot two standardized series in the same graph. Figure 4 shows in most of the time period, google trends variables have high correlations with the sales data.

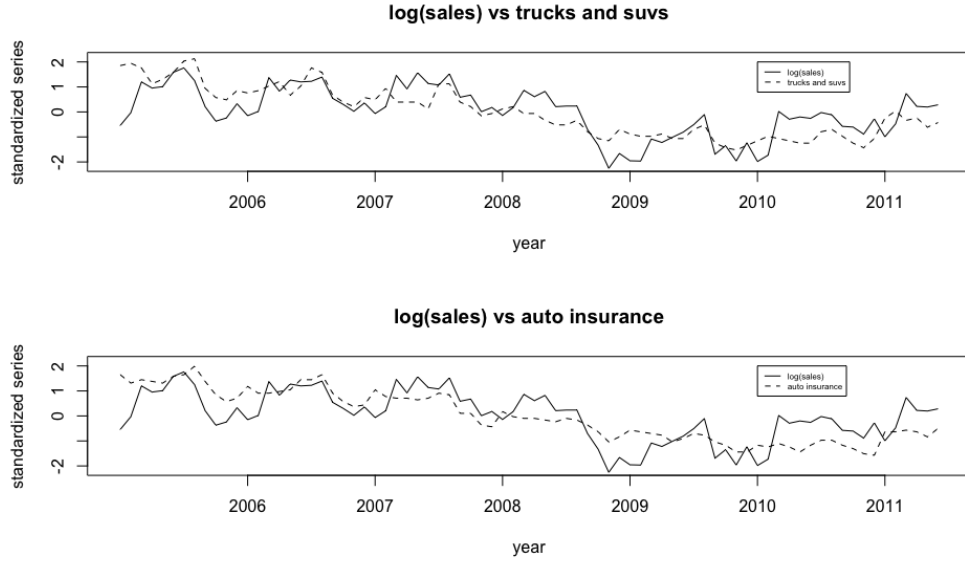


Figure 4: temporal aspects relationship

Therefore, the google trends model becomes

$$y_t = b_0 + b_1 y_{t-1} + b_{12} y_{t-12} + a_1 g_{1,t} + a_2 g_{2,t} + e_t \quad (\text{google trends model})$$

The  $R^2$  of the google trends model is 0.7821, which suggests the google trends variables significantly improve in-sample fit when added to this regression. Moreover, the diagnosis of the google trends model is based on the graphs of residuals.

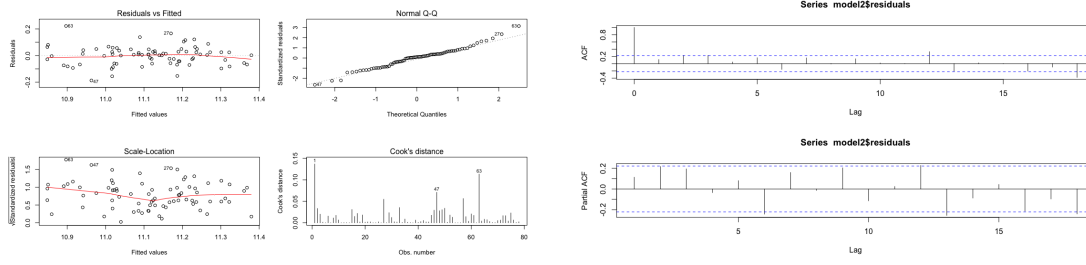


Figure 5: Diagnosis of the google trends model

It seems that the residual is the white-noise series with normal distribution, indicating the google trends model is a suitable model in this case. The further question of interest is whether the trends variables improve out-of-sample forecasting. To check this, we use a rolling window forecast. We set the length of each rolling window to be two-year and make an one-step ahead prediction. Then we compute the MAE for each model and make a comparison. The results are shown in figure 6. The overall MAE decreases from 6.38% to 5.71%. However, if we just examine the MAE especially in the recession period, it decreases from 7.61% to 5.19%.

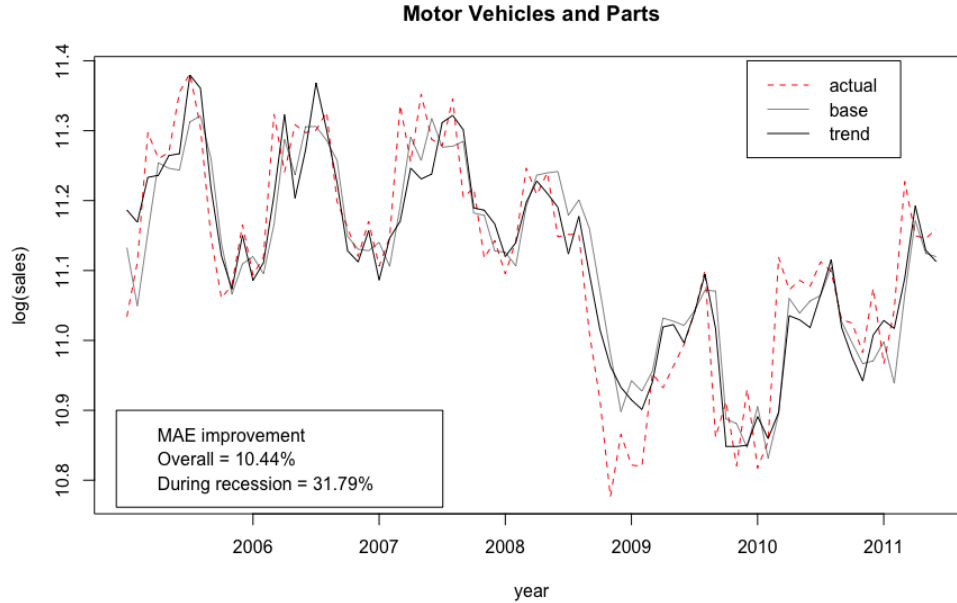


Figure 6: Motor Vehicles and Parts

So far we notice google trends can indeed help us predict the present. The last question of interest in this section is how well the google trends model predict the present. To investigate this question, we try to leave out the end of the actual data set and use Google trends to make one-step ahead prediction for the period 2005-06-01 to 2011-06-01. Figure 7 shows the prediction is quite well in most of the period. It may give some indication of the peak points and change points. The worst predicted period is the end of the 2009 year. Figure 4 may explain part of the reason. In that period the relationship between google trends and sales is not really strong probably due to some unknown events. But overall the google trends predict the present sales quite well.

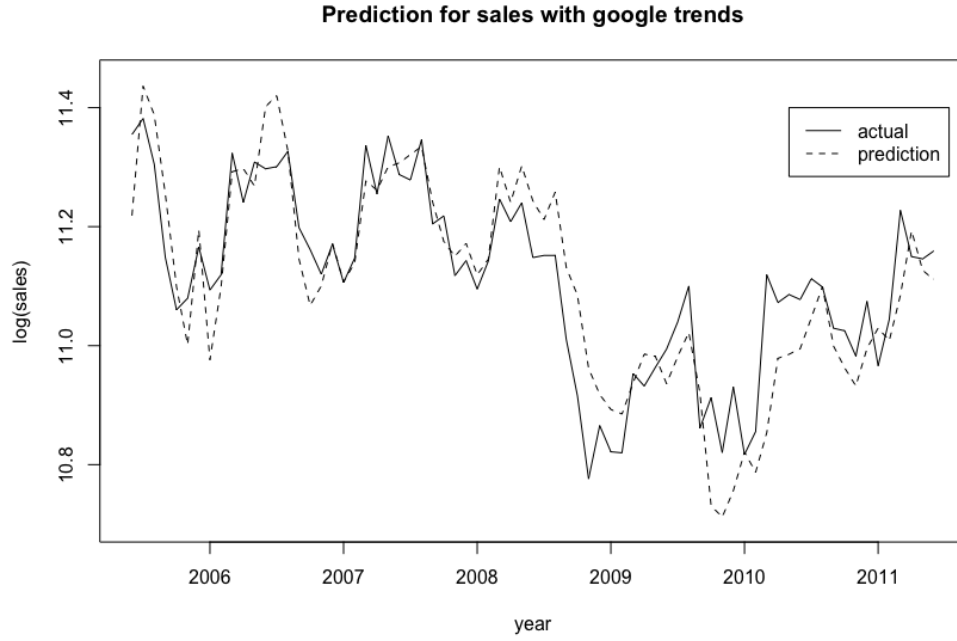


Figure 7: Prediction for sales with google trends

## 5 CDC Flu Data Example

Each week, the CDC estimates the national percentage of visits to healthcare providers for influenza like illnesses. This measure (denoted ILINet) is used to assess questions of flu season start, end, intensity, peak times, etc... Understandably, having the most up to date information on ILINet is a priority; but, given the complexity of measuring ILINet, reporting delays are inevitable. However, perhaps using Google Trends data for the previous week, the most up-to-date value of ILINet can be predicted. We provide a graph of the relevant data in the figure 8.

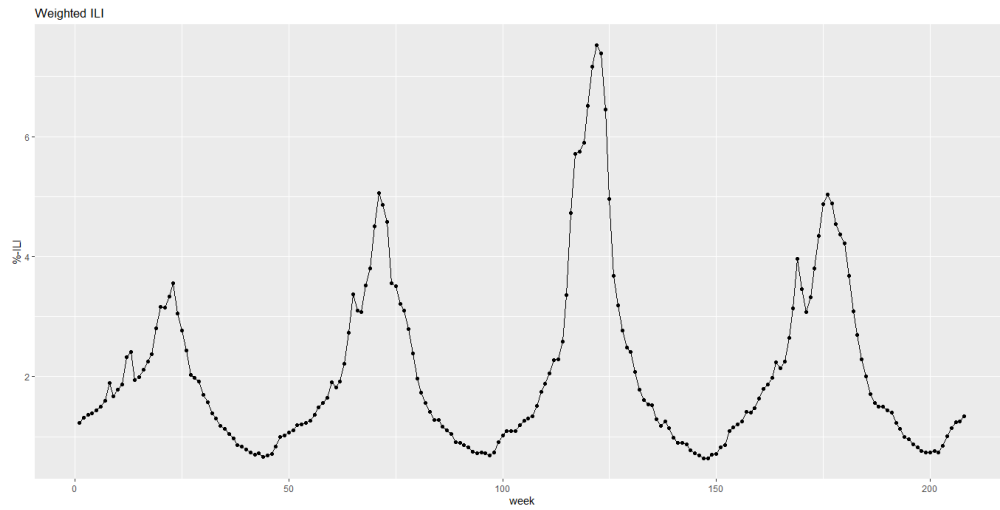


Figure 8: FLU Data

Inspection of the graph suggests that a seasonal AR model might be appropriate. The ACF and PACF plots (shown in the figure 9) further indicate that a seasonal AR(2) model should be adequate. Furthermore, running a Box-Cox procedure (not shown) suggests heavily that an inverse transformation is needed. We thus fit the seasonal AR(2) model  $x_t = b_0 + b_1x_{t-1} + b_2x_{t-2} + b_{52}x_{t-52} + e_t$  to the inverse of the ILINet data. The model provides a very good fit, with an adjusted  $R^2$  of 0.9833.

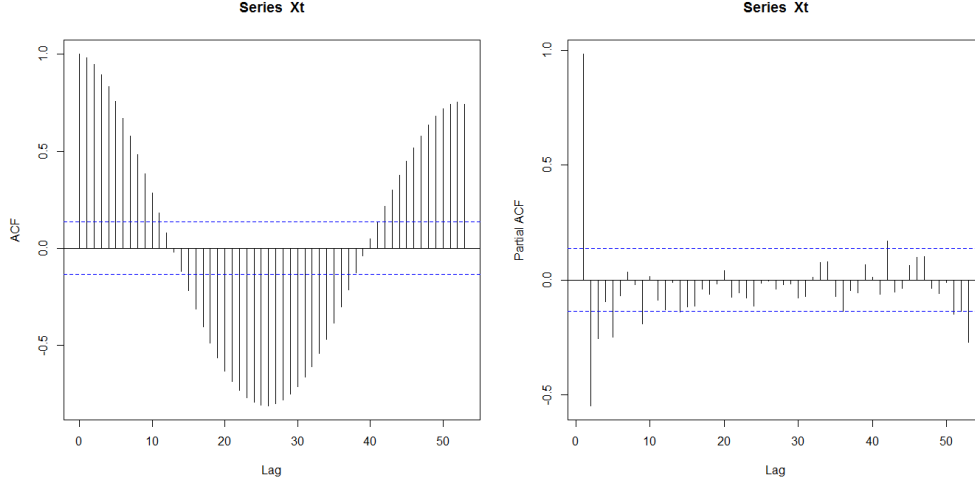


Figure 9: ACF and PACF plot

There are many searches that could be indicative of flu activity, such as “flu season,” “flu shots,” “flu symptoms,” and many more. We considered a best subsets routine for models including “flu season,” “flu shots,” “flu symptoms,” and “flu” search query data (for the U.S.) as exogenous predictors, as well as the lag values assumed in the base model. The routine suggested that the model utilizing all but the “flu shots” data gives the best possible model by AIC and adjusted  $R^2$  criterion. See the code supplied in the appendix . We therefore fit a seasonal AR(2) model with these three Google Search terms as exogenous predictors. The model equation is  $x_t = b_0 + b_1x_{t-1} + b_2x_{t-2} + b_{52}x_{t-52} + a_1g_{1,t} + a_2g_{2,t} + a_3g_{3,t}e_t$  (The BIC criterion suggests that only the “flu shots” data should be used). The Adjusted  $R^2$  for this new model is 0.9842.

We note that there is only a very small improvement to in-sample fit when these predictors are added to the base model. Perhaps a significant improvement in predictive ability will be seen if we consider out of sample forecasts instead. To check this, we use a rolling window forecast for each model, beginning using the window  $t=1$  through  $t=104$  (the first two seasons) and forecasting the ILINet value at  $t=105$ , then shifting the window over by 1 and repeating the procedure of forecasting until we forecast the last value of ILINet available. We then compute the mean absolute error (MAE) for each model’s set of forecasts. We find the MAE for the base model is 3.996544%, while that for the model utilizing google trends data is 4.254782%, representing an overall 6.461% *reduction* in fit. However, closer examination of the fit reveals a potentially salvageable result.

While overall predictive ability is certainly an interest, there are more specific aspects of this flu data that are of interest to public health officials. In particular, the values of ILINet around peak season (when ILINet is at its highest) are of special interest to public health officials, for these values relate to the most severe week of flu activity, and how severe it will be during this

time, which has consequences for advisory warnings issued to hospitals and citizens. We therefore have reason to consider more closely the periods of time around peak season for each season. We compute the MAE for each model using forecasts five weeks before, five weeks after, and at the peak week for the final 3 seasons captured in the data at hand. A rather encouraging picture then emerges. Two of the three seasons saw an improvement in fit, with one season experiencing a drastic improvement. More specifically, a 12.03294% improvement in MAE was seen in the first forecast season, a staggering 63.71317% improvement in MAE was seen in the second forecast season, and an unfortunate 13.26341% reduction in fit was seen in the third and final forecast season. (The model selected under the BIC criterion actually has a small improvement to MAE of approximately 1.8%. Additionally, it shows an marked improvement around season peaks, though these improvements are less pronounced.) We summarize these results in a graph of in sample fit (figure 10). Note the improvemnt of the Google Trends model around peak week for the first two seasons.

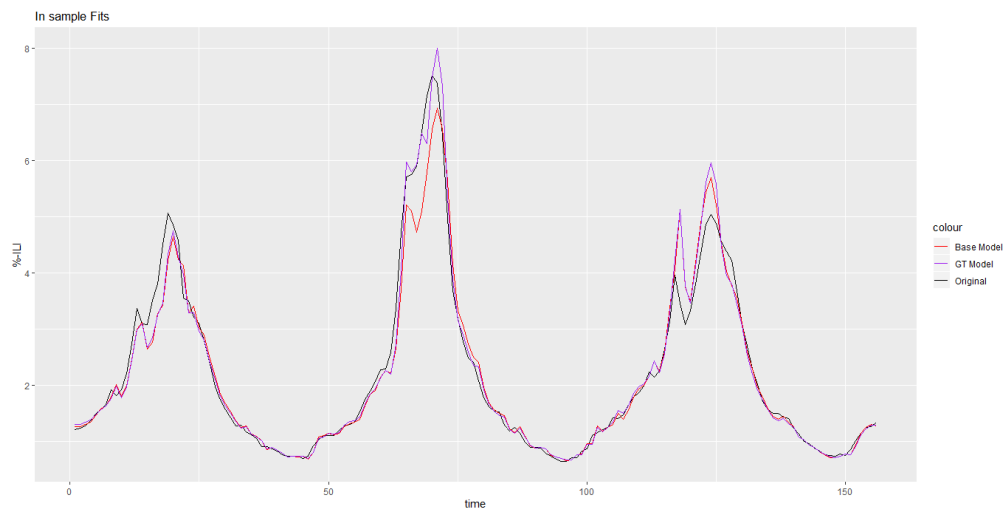


Figure 10: Final Comparison

It might be wondered what the results of trying to predict the current ILINet value using only google trends data would be. To this end, we consider building a model using search data for the query “flu.” Since Google Trends broad matches search queries, the search for “flu” should encode most information related to flu-based searches. We run a best subset routine considering current values for the “flu” query search, as well as lag values of 1, 2 and 52. (See code provided in the appendix for details.) Using the AIC criterion, the model with all predictors is considered optimal. This model has a markedly decreased sample fit (adjusted  $R^2 = 0.4768$ ). We then consider the same rolling window forecast considered above. The results are less than promising, with a MAE calculated as 0.254087, representing a staggering 535.7669% reduction in fit compared to the original baseline seasonal AR(2) model. We would therefore advise against using Google Trends data alone as a predictor in this case.

## 6 Conclusion

Two examples of how Google Trends data may be used to predict the present have been discussed. In one, we considered how Google Trends data might be used to forecast current values of economic indicators, such as car sales; in the other, we considered how Google Trends might be



used to forecast flu season. In the automobiles example, we found that both in sample fit, out of sample forecasting ability were improved by the inclusion of google trends data. In the flu season example, neither in sample fit, nor out of sample predictive ability were greatly improved. However, inclusion of Google Trends data did provide a significant improvement in predictive ability for certain specific aspects of the data, namely around each season's peak week. Similar improvements for specific aspects of the automobiles data, namely during the December 2007 through June 2009 recession, were also seen in the automobiles sales example when Google Trends data were considered.

These results demonstrate the potential usefulness of Google Trends in forecasting applications. If a base model doesn't have an extremely high in sample fit, adding Google search data could potentially improve that fit, and subsequent forecasting accuracy. Furthermore, even if in sample fit doesn't greatly benefit from the addition of Google Trends data, forecasting accuracy for specific aspects of the data of interest could still be improved.

We close with a suggestion for further research. It has been noted that combining sets of forecasts made with different models can be potentially more accurate (in terms of MSE) than individual forecasts made with a given model.[2] It is possible, then, that composite forecasting models can be built using several models built with google trends data to improve forecasting accuracy for different aspects of the data of interest (such as season peak for the flu data above) that yield overall increased accuracy for each aspect of the data considered. Thus, even greater accuracy can potentially be achieved using Google Trends data than seen in the examples presented here. We therefore encourage research be devoted to ascertaining whether this composite forecasting method could work as intuition would seem to suggest it might.

## References

- [1] Choi, H., & Varian, H. (2012). Predicting the present with Google Trends. *Economic Record*, 88, 2-9.
- [2] Bates, J., & Granger, C. (1969). The Combination of Forecasts. *OR*, 20(4), 451-468.

## Appendix A - Code for automobiles example

```
1 # base model
2 data1 = read.csv("/Users/lixuezhen/Desktop/237\ project/vehicles.csv",header = T,skip = 6)
3 head(data1)
4 y = log(data1$Value)
5 plot.ts(y[13:90],xlab = "year", xaxt = "n",main = "Motor Vehicles and Parts Dealers",ylab =
  "log(sales)")
6 axis(side = 1, at = c(13,25,37,49,61,73),labels = c(2006,2007,2008,2009,2010,2011))
7 acf(y[13:90],lag.max = 30,main = "ACF plot")
8 model1 = lm(y[13:90]~y[12:89]+y[1:78])
9 summary(model1)
10 library(leaps)
11 model1_2 = lm(y[13:90]~y[12:89]+y[11:88]+y[10:87]+y[9:86]+y[8:85]+y[1:78])
12 summary(model1_2)
13 fit_all = regsubsets(y[13:90]~y[12:89]+y[11:88]+y[10:87]+y[9:86]+y[8:85]+y[1:78],data = as.
  data.frame(y),nbest =1, method="exhaustive")
14 fit_all.display = as.data.frame(summary(fit_all)$outmat)
15 fit_all.display = cbind(fit_all.display,round(summary(fit_all)$rsq,4),round(summary(fit_all)
  $adjr2,4),round(summary(fit_all)$cp,4),round(summary(fit_all)$bic,4))
16 varnames = c("lag1","lag2","lag3","lag4","lag5","lag12")
17 names(fit_all.display)=c(varnames,"R^2","adj R^2","Cp","BIC")
18 fit_all.display
19
20 #diagnosis
21 par(mfrow = c(2,2))
22 plot(model1,which = 1:4)
23 par(mfrow = c(1,1))
24 acf(model1$residuals)
25 pacf(model1$residuals)
26
27
28 # google trends model
29 trucks_suvs = read.csv("/Users/lixuezhen/Desktop/237\ project/trucks_suvs.csv",header = T,
  skip = 1)
30 auto_insurance = read.csv("/Users/lixuezhen/Desktop/237\ project/auto_insurance.csv",header
  = T,skip = 1)
31 head(trucks_suvs)
32 head(auto_insurance)
33 suvs = trucks_suvs$Geo..United.States
34 insurance = auto_insurance$Geo..United.States
35 par(mfrow = c(2,1))
36 plot.ts((y[13:90]-mean(y[13:90]))/sd(y[13:90]),xlab = "year", ylim = c(-2.2,2.2),xaxt = "n",
  main = "log(sales) vs trucks and suvs",ylab = "standardized series")
37 axis(side = 1, at = c(13,25,37,49,61,73),labels = c(2006,2007,2008,2009,2010,2011))
38 lines((suvs[13:90]-mean(suvs[13:90]))/sd(suvs[13:90]),lty = 2)
39 legend(61,2,legend = c("log(sales)","trucks and suvs"),lty = c(1,2),cex = 0.5)
40 plot.ts((y[13:90]-mean(y[13:90]))/sd(y[13:90]),xlab = "year", ylim = c(-2.2,2.2),xaxt = "n",
  main = "log(sales) vs auto insurance",ylab = "standardized series")
41 axis(side = 1, at = c(13,25,37,49,61,73),labels = c(2006,2007,2008,2009,2010,2011))
42 lines((insurance[13:90]-mean(insurance[13:90]))/sd(insurance[13:90]),lty = 2)
43 legend(61,2,legend = c("log(sales)","auto insurance"),lty = c(1,2),cex = 0.5)
44 par(mfrow = c(1,1))
45 model2 = lm(y[13:90]~y[12:89]+y[1:78]+suvs[13:90]+insurance[13:90])
46 summary(model2)
47
48 #diagnosis
49 par(mfrow = c(2,2))
50 plot(model2,which = 1:4)
51 par(mfrow = c(1,1))
52 par(mfrow = c(2,1))
53 acf(model2$residuals)
54 pacf(model2$residuals)
55 par(mfrow = c(1,1))
56
57 # final plots
```

```

58 plot.ts(y[13:90],ylab = "log(sales)",xlab = "year",xaxt = "n",col = "red",lty = 2,main = "
    Motor Vehicles and Parts")
59 axis(side = 1, at = c(13,25,37,49,61,73),labels = c(2006,2007,2008,2009,2010,2011))
60 lines(model1$fitted.values,lwd = 0.5,lty = 1)
61 lines(model2$fitted.values,lwd = 1,lty = 1)
62 legend(60,11.4,legend = c("actual","base","trend"),col = c("red","black","black"),lwd = c
    (1,0.5,1),lty = c(2,1,1))
63 legend(0,10.9,legend = c("MAE improvement","Overall = 10.44%","During recession = 31.79%"))
64
65 # MAE comparison
66 error1 = NULL
67 for (i in 31:90){
68     fit = lm(y[(i-18):(i-1)]~y[(i-19):(i-2)]+y[(i-30):(i-13)])
69     error1 = c(error1,sum(fit$coefficients * c(1,y[i-1],y[i-12])) - y[i])
70 }
71 mae1 = sum(abs(error1))/length(error1)
72 mae1
73 error2 = NULL
74 for (i in 31:90){
75     fit = lm(y[(i-18):(i-1)]~y[(i-19):(i-2)]+y[(i-30):(i-13)]+suvs[(i-18):(i-1)]+insurance[(i
        -18):(i-1)])
76     error2 = c(error2,sum(fit$coefficients * c(1,y[i-1],y[i-12],suvs[i],insurance[i])) - y[i])
77 }
78
79 mae2 = sum(abs(error2))/length(error2)
80 mae2
81 (mae1-mae2)/mae1 # overall improvement
82
83 mae1 = sum(abs(error1[18:35]))/length(error1[18:35])
84 mae2 = sum(abs(error2[18:35]))/length(error2[18:35])
85 (mae1-mae2)/mae1 # recession period improvement
86
87 # leave out the the end of the actual data set
88 predict_value = NULL
89 for (i in 18:90){
90     fit = lm(y[13:(i-1)]~y[12:(i-2)]+y[1:(i-13)]+suvs[13:(i-1)]+insurance[13:(i-1)])
91     predict_value = c(predict_value,sum(fit$coefficients * c(1,y[i-1],y[i-12],suvs[i],insurance[
        i]))))
92 }
93 plot.ts(y[18:90],xlab = "year",ylab= "log(sales)",xaxt = "n",main = "Prediction for sales
    with google trends",ylim = c(10.7,11.45))
94 axis(side = 1, at = c(8,20,32,44,56,68),labels = c(2006,2007,2008,2009,2010,2011))
95 lines(predict_value,lty = 2)
96 legend(60,11.4,legend = c("actual","prediction"),lty = c(1,2))
97
98 mean(abs(predict_value-y[18:90]))

```

## Appendix B - Code for flu data example

```

1 library(ggplot2)
2 library(leaps)
3 library(MASS)
4 setwd("C:/Users/Kenneth/Downloads/Flu (15.16-18.19)")
5 flu<-read.csv("ILINetNational.csv", header=TRUE)
6
7 n<-length(flu$X.WEIGHTED.ILI)
8 p1<-ggplot(data=flu, aes(x=1:n,y=X.WEIGHTED.ILI))+geom_line()+geom_point()
9 p1+labs(title="Weighted ILI", y="%-ILI", x="week")
10 x11()
11 Xt<-1/flu$X.WEIGHTED.ILI
12 p2<-ggplot(data=as.data.frame(Xt), aes(x=1:n,y=Xt))+geom_line()+geom_point()
13 p2+labs(title="Inverse ILI", y="1/%-ILI", x="week")
14 model1<-lm(flu$X.WEIGHTED.ILI~L1+L2+L52)
15 boxcox(model1)##This suggests an inverse transformation, which we use below.
16
17 ##ACF and PACF plots; suggestive of an AR(2) model.

```

```

18 par(mfrow=c(1,2))
19 acf(Xt,lag=53)
20 pacf(Xt,lag=53)
21
22 ##Model fit with Regression approach seen in the paper:
23 L1<-c(NA,Xt[1:(NROW(Xt)-1)])
24 L2<-c(NA,NA,Xt[1:(NROW(Xt)-2)])
25 L52<-c(rep(NA,52),Xt[1:(NROW(Xt)-52)])
26
27 model1<-lm(Xt~L1+L2+L52)
28 summary(model1)
29 par(mfrow=c(2,2))
30 plot(model1, which=1)
31 plot(model1, which=2)
32 plot(model1, which=3)
33 plot(model1, which=5)
34
35 np<-length(model1$fitted)
36 ggplot()+geom_line(data=as.data.frame(model1$fitted), aes(x=1:np,y=model1$fitted), color='
red')+
37 geom_line(data=as.data.frame(Xt[(n+1-np):n]), aes(x=1:np,y=Xt[(n+1-np):n]), color='black')
38
39 ##Once a base model has been chosen, we can use GT data as additional predictors:
40 fluseason<-read.csv("fluseasonrends.csv", header=TRUE)
41 flushot<-read.csv("flushottrends.csv", header=TRUE)
42 flusymptoms<-read.csv("flusymptomstends.csv", header=TRUE)
43 flutrends<-read.csv("flutrends.csv", header=TRUE)
44
45 DATA<-as.data.frame(cbind(Xt,L1,L2,L52,fluseason,flushot,flusymptoms,flutrends))
46 best<-regsubsets(Xt~L1+L2+L52+flu.season+flu.shot+flu.symptoms+flu, data=DATA)
47 ss1<-summary(best)
48 n1<-length(Xt)
49 p1<-as.integer(rownames(ss1$which))+1
50 aicp1<-n1*log(ss1$rss/n1)+2*p1
51 object1<-cbind(ss1$which, ss1$rss, ss1$rsq, ss1$adjr2, ss1$cp, aicp1, ss1$bic)
52 colnames(object1)<-c(colnames(ss1$which), "sse", "Rsq", "adjRsq", "Cp", "AIC", "BIC")
53 object1
54
55 model2<-lm(Xt~L1+L2+L52+fluseason$flu.season+flusymptoms$flu.symptoms+flutrends$flu)
56 ##This appears to be the best model using utilizing the GT data AIC Criterion
57 summary(model2)
58 ##Rolling forecasts:
59 ##Model1 Forecasts; N=n-m subsamples
60 m<-104
61 error<-NULL
62 for(i in 1:m){
63 XR<-Xt[i:(i+m-1)]
64 lag1<-c(NA,XR[1:(NROW(XR)-1)])
65 lag2<-c(NA,NA,XR[1:(NROW(XR)-2)])
66 lag52<-c(rep(NA,52),XR[1:(NROW(XR)-52)])
67 fit<-lm(XR~lag1+lag2+lag52)
68 error[i]<-(sum(fit$coefficients*c(1,Xt[i+m-1],Xt[i+m-2],Xt[i+m-52])))-Xt[i+m]
69 }
70 MAE1<-sum(abs(error))/length(error)
71 MAE1
72
73 ##Model2 rolling forecasts
74 error2<-NULL
75 for(i in 1:m){
76 XR<-Xt[i:(i+m-1)]
77 lag1<-c(NA,XR[1:(NROW(XR)-1)])
78 lag2<-c(NA,NA,XR[1:(NROW(XR)-2)])
79 lag52<-c(rep(NA,52),XR[1:(NROW(XR)-52)])
80 trend1<-fluseason$flu.season[i:(i+m-1)]
81 trend2<-flusymptoms$flu.symptoms[i:(i+m-1)]
82 trend3<-flutrends$flu[i:(i+m-1)]
83 fit<-lm(XR~lag1+lag2+lag52+trend1+trend2+trend3, data=NULL)

```

```

84 error2[i]<-(sum(fit$coefficients*c(1,Xt[i+m-1],Xt[i+m-2],Xt[i+m-52],fluseason$flu.season[i+m
    ],flusymptoms$flu.symptoms[i+m],flutrends$flu[i+m]))) -Xt[i+m]
85 }
86 MAE2<-sum(abs(error2))/length(error2)
87 MAE2
88
89 (MAE1-MAE2)/MAE1
90
91 ##MAE Comparison around peaks
92 m1<-max(flu$X.WEIGHTED.ILI[53:104])
93 m2<-max(flu$X.WEIGHTED.ILI[105:156])
94 m3<-max(flu$X.WEIGHTED.ILI[157:208])
95
96 flu$X.WEIGHTED.ILI==m1
97 flu$X.WEIGHTED.ILI==m2
98 flu$X.WEIGHTED.ILI==m3
99 ##The last three peaks occur at points 71, 122, and 176
100
101 error1b<-NULL
102 error1g<-NULL
103 error2b<-NULL
104 error2g<-NULL
105 error3b<-NULL
106 error3g<-NULL
107 for(i in -5:5){
108 error1b[i+6]<-model1$fitted[i+71-52]-Xt[i+71]
109 error1g[i+6]<-model2$fitted[i+71-52]-Xt[i+71]
110 error2b[i+6]<-model1$fitted[i+122-52]-Xt[i+122]
111 error2g[i+6]<-model2$fitted[i+122-52]-Xt[i+122]
112 error3b[i+6]<-model1$fitted[i+176-52]-Xt[i+176]
113 error3g[i+6]<-model2$fitted[i+176-52]-Xt[i+176]
114 }
115 MAEb1<-sum(abs(error1b))/length(error1b)
116 MAEg1<-sum(abs(error1g))/length(error1g)
117 MAEb2<-sum(abs(error2b))/length(error2b)
118 MAEg2<-sum(abs(error2g))/length(error2g)
119 MAEb3<-sum(abs(error3b))/length(error2b)
120 MAEg3<-sum(abs(error3g))/length(error3g)
121
122 (MAEb1-MAEg1)/MAEb1
123 (MAEb2-MAEg2)/MAEb2
124 (MAEb3-MAEg3)/MAEb3
125
126 np<-length(model1$fitted)
127 ##Model in-sample fit comparisons
128 plot<-ggplot()+geom_line(data=as.data.frame(Xt[(n+1-np):n]), aes(x=1:np,y=1/Xt[(n+1-np):n],
    color="Original"))+
129 geom_line(data=as.data.frame(model1$fitted), aes(x=1:np,y=1/model1$fitted,color="Base Model"
    ))+
130 geom_line(data=as.data.frame(model2$fitted), aes(x=1:np,y=1/model2$fitted, color="GT Model")
    )
131 plot+labs(title="In sample Fits", y="%-ILI", x="time")+
132 scale_color_manual(values=c("red", "purple", "black"))
133
134
135 ##Predict flu using flu shot GT data only:
136 St<-flutrends$flu
137 ggplot()+geom_line(data=flu, aes(x=1:n,y=X.WEIGHTED.ILI), color='black')+
138 geom_line(data=as.data.frame(St), aes(x=1:n,y=St), color='red')
139
140 g11<-c(NA,St[1:(NROW(St)-1)])
141 g12<-c(NA,NA,St[1:(NROW(St)-2)])
142 g152<-c(rep(NA,52),St[1:(NROW(St)-52)])
143 subset<-regsubsets(Xt~St+g11+g12+g152,data=NULL)
144 ss2<-summary(subset)
145 n2<-length(St)
146 p2<-as.integer(rownames(ss2$which))+1
147 aicp2<-n*log(ss2$rss/n2)+2*p2

```

```

148 object2<-cbind(ss2$which, ss2$rss, ss2$rsq, ss2$adjr2, ss2$cp, aicp2, ss2$bic)
149 colnames(object2)<-c("int", "lag0", "lag1", "lag2", "lag52", "sse", "Rsq", "adjRsq", "Cp", "AIC"
150 , "BIC")
151 object2
152 ##Model with all predictors apperas to be the best (by AIC criterion)
153 ##GT model rolling forecasts
154 error3<-NULL
155 for(i in 1:m){
156   GTR<-St[i:(i+m-1)]
157   XR<-Xt[i:(i+m-1)]
158   lag1<-c(NA, GTR[1:(NROW(GTR)-1)])
159   lag2<-c(NA, NA, GTR[1:(NROW(GTR)-2)])
160   lag52<-c(rep(NA, 52), GTR[1:(NROW(GTR)-52)])
161   fit<-lm(XR~GTR+lag1+lag2+lag52, data=NULL)
162   error3[i]<-(sum(fit$coefficients*c(1, St[i+m], St[i+m-1], St[i+m-2], St[i+m-52]))) -Xt[i+m]
163 }
164 MAE3<-sum(abs(error3))/length(error3)
165 MAE3
166
167 (MAE1-MAE3)/MAE1 ##Note the extreme reduction in fit, reflected in the graph below.
168
169 model3<-lm(Xt~St+gl1+gl2+gl52, data=NULL)
170 np2<-length(model3$fitted)
171 ##Model in sample fit comparisons
172 plot2<-ggplot()+geom_line(data=as.data.frame(Xt[(n+1-np2):n]), aes(x=1:np2, y=1/Xt[(n+1-np2):
173 n], color="Original"))+
174 geom_line(data=as.data.frame(model1$fitted), aes(x=1:np, y=1/model1$fitted, color="Base Model"
175 ))+
176 geom_line(data=as.data.frame(model3$fitted), aes(x=1:np2, y=1/model3$fitted, color="GT Pure
177 Model"))
178
179 plot2+labs(title="In sample Fits", y="%-ILI", x="time")+
180 scale_color_manual(values=c("red", "purple", "black"))
181
182 ##-----
183 ##The following code may be used to run the analysis discussed in the Flu example
184 ##using the google trends model selected by the BIC criterion.
185 model2<-lm(Xt~L1+L2+L52+flushot$flu.shot)
186 ##This appears to be the best model using utilizing the GT data BIC criterion
187 summary(model2)
188 ##Rolling forecasts:
189 ##Model1 Forecasts; N=n-m subsamples
190 m<-104
191 error<-NULL
192 for(i in 1:m){
193   XR<-Xt[i:(i+m-1)]
194   lag1<-c(NA, XR[1:(NROW(XR)-1)])
195   lag2<-c(NA, NA, XR[1:(NROW(XR)-2)])
196   lag52<-c(rep(NA, 52), XR[1:(NROW(XR)-52)])
197   fit<-lm(XR~lag1+lag2+lag52)
198   error[i]<-(sum(fit$coefficients*c(1, Xt[i+m-1], Xt[i+m-2], Xt[i+m-52]))) -Xt[i+m]
199 }
200 MAE1<-sum(abs(error))/length(error)
201 MAE1
202
203 ##Model2 rolling forecasts
204 error2<-NULL
205 for(i in 1:m){
206   XR<-Xt[i:(i+m-1)]
207   lag1<-c(NA, XR[1:(NROW(XR)-1)])
208   lag2<-c(NA, NA, XR[1:(NROW(XR)-2)])
209   lag52<-c(rep(NA, 52), XR[1:(NROW(XR)-52)])
210   trend1<-flushot$flu.shot[i:(i+m-1)]
211   fit<-lm(XR~lag1+lag2+lag52+trend1, data=NULL)
212   error2[i]<-(sum(fit$coefficients*c(1, Xt[i+m-1], Xt[i+m-2], Xt[i+m-52], flushot$flu.shot[i+m])))
213     -Xt[i+m]}
214 MAE2<-sum(abs(error2))/length(error2)

```

```

211 MAE2
212
213 (MAE1-MAE2)/MAE1
214
215 ##MAE Comparison around peaks
216 m1<-max(flu$X.WEIGHTED.ILI[53:104])
217 m2<-max(flu$X.WEIGHTED.ILI[105:156])
218 m3<-max(flu$X.WEIGHTED.ILI[157:208])
219
220 flu$X.WEIGHTED.ILI==m1
221 flu$X.WEIGHTED.ILI==m2
222 flu$X.WEIGHTED.ILI==m3
223 ##The last three peaks occur at points 71, 122, and 176
224
225 error1b<-NULL
226 error1g<-NULL
227 error2b<-NULL
228 error2g<-NULL
229 error3b<-NULL
230 error3g<-NULL
231 for(i in -5:5){
232 error1b[i+6]<-model1$fitted[i+71-52]-Xt[i+71]
233 error1g[i+6]<-model2$fitted[i+71-52]-Xt[i+71]
234 error2b[i+6]<-model1$fitted[i+122-52]-Xt[i+122]
235 error2g[i+6]<-model2$fitted[i+122-52]-Xt[i+122]
236 error3b[i+6]<-model1$fitted[i+176-52]-Xt[i+176]
237 error3g[i+6]<-model2$fitted[i+176-52]-Xt[i+176]
238 }
239 MAEb1<-sum(abs(error1b))/length(error1b)
240 MAEg1<-sum(abs(error1g))/length(error1g)
241 MAEb2<-sum(abs(error2b))/length(error2b)
242 MAEg2<-sum(abs(error2g))/length(error2g)
243 MAEb3<-sum(abs(error3b))/length(error2b)
244 MAEg3<-sum(abs(error3g))/length(error3g)
245
246 (MAEb1-MAEg1)/MAEb1
247 (MAEb2-MAEg2)/MAEb2
248 (MAEb3-MAEg3)/MAEb3

```