

Predicting Flu (Linear Regression)

Saeid Abolfazli

May 2, 2016

Flu epidemics constitute a major public health concern causing respiratory illnesses, hospitalizations, and deaths. According to the National Vital Statistics Reports published in October 2012, influenza ranked as the eighth leading cause of death in 2011 in the United States. Each year, 250,000 to 500,000 deaths are attributed to influenza related diseases throughout the world.

The U.S. Centers for Disease Control and Prevention (CDC) and the European Influenza Surveillance Scheme (EISS) detect influenza activity through virologic and clinical data, including Influenza-like Illness (ILI) physician visits. Reporting national and regional data, however, are published with a 1-2 week lag.

The Google Flu Trends project was initiated to see if faster reporting can be made possible by considering flu-related online search queries – data that is available almost immediately.

Problem 1.1 - Understanding the Data

We would like to estimate influenza-like illness (ILI) activity using Google web search logs. Fortunately, one can easily access this data online:

ILI Data - The CDC publishes on its website the official regional and state-level percentage of patient visits to healthcare providers for ILI purposes on a weekly basis.

Google Search Queries - Google Trends allows public retrieval of weekly counts for every query searched by users around the world. For each location, the counts are normalized by dividing the count for each query in a particular week by the total number of online search queries submitted in that location during the week. Then, the values are adjusted to be between 0 and 1.

The csv file `FluTrain.csv` aggregates this data from January 1, 2004 until December 31, 2011 as follows:

- **“Week”** - The range of dates represented by this observation, in year/month/day format.
- **“ILI”** - This column lists the percentage of ILI-related physician visits for the corresponding week.
- **“Queries”** - This column lists the fraction of queries that are ILI-related for the corresponding week, adjusted to be between 0 and 1 (higher values correspond to more ILI-related search queries).

Before applying analytics tools on the training set, we first need to understand the data at hand. Load “`FluTrain.csv`” into a data frame called `FluTrain`.

```
FluTrain <- read.csv("data/FluTrain.csv")
```

Looking at the time period 2004-2011, which week corresponds to the highest percentage of ILI-related physician visits? Select the day of the month corresponding to the start of this week.

```
FluTrain[which.max(FluTrain$ILI),]
```

```
##                Week      ILI Queries
## 303 2009-10-18 - 2009-10-24 7.618892      1
```

Which week corresponds to the highest percentage of ILI-related query fraction?

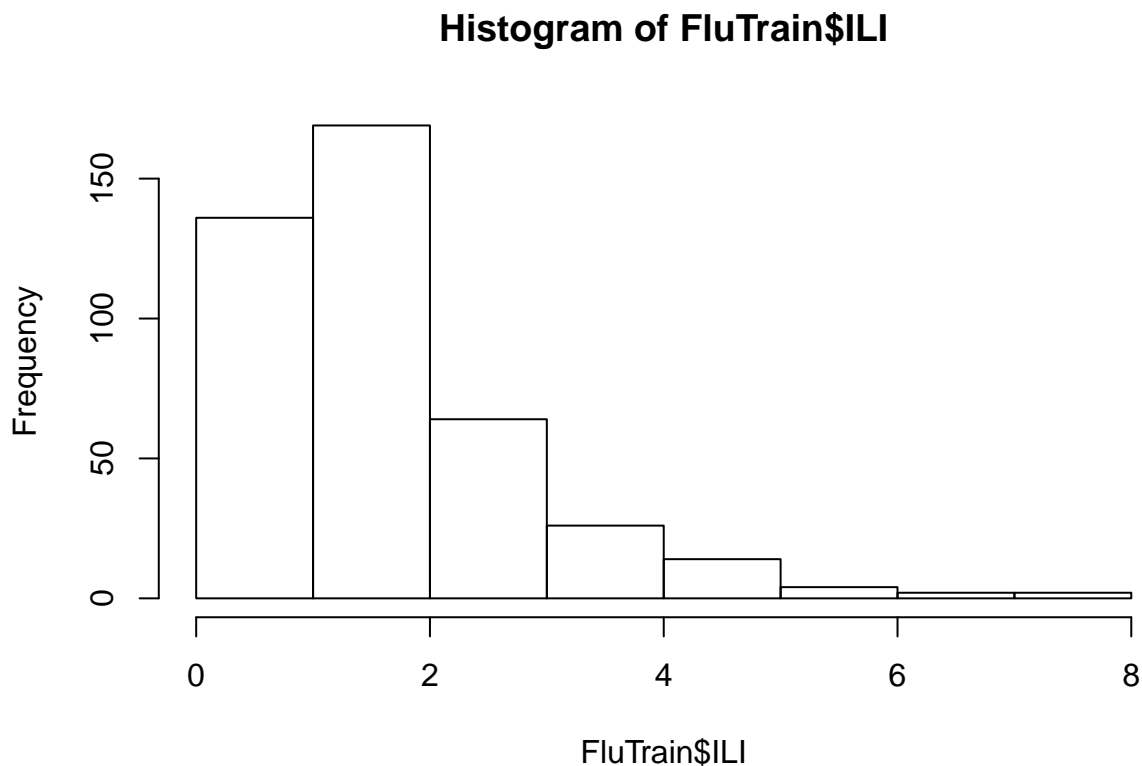
```
FluTrain[which.max(FluTrain$Queries),]
```

```
##                Week      ILI Queries
## 303 2009-10-18 - 2009-10-24 7.618892      1
```

Problem 1.2 - Understanding the Data

Let us now understand the data at an aggregate level. Plot the histogram of the dependent variable, ILI. What best describes the distribution of values of ILI?

```
hist(FluTrain$ILI)
```



1. Most of the ILI values are small, with a relatively small number of much larger values (in statistics, this sort of data is called “skew right”).
2. The ILI values are balanced, with equal numbers of unusually large and unusually small values.
3. Most of the ILI values are large, with a relatively small number of much smaller values (in statistics, this sort of data is called “skew left”).

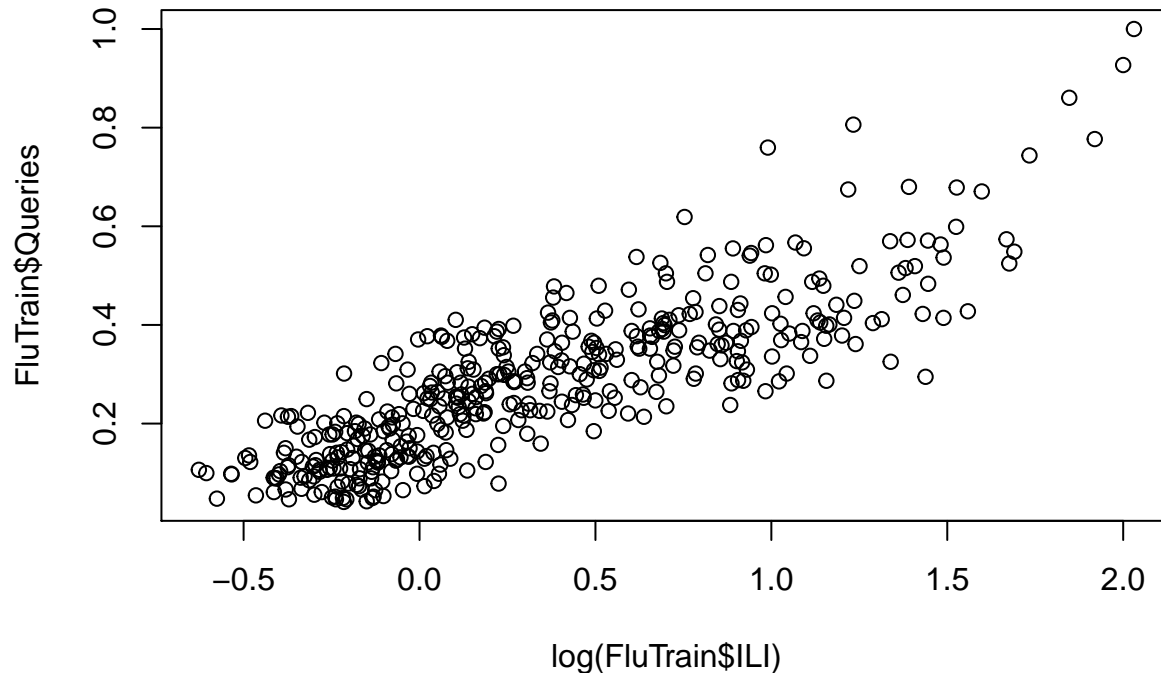
Problem 1.3 - Understanding the Data

When handling a skewed dependent variable, it is often useful to predict the logarithm of the dependent variable instead of the dependent variable itself – this prevents the small number of unusually large or small

observations from having an undue influence on the sum of squared errors of predictive models. In this problem, we will predict the natural log of the ILI variable, which can be computed in R using the `log()` function.

Plot the natural logarithm of ILI versus Queries. What does the plot suggest?.

```
plot(log(FluTrain$ILI), FluTrain$Queries)
```



1. There is a negative, linear relationship between `log(ILI)` and Queries.
2. There is no apparent linear relationship between `log(ILI)` and Queries.
3. There is a positive, linear relationship between `log(ILI)` and Queries.

Answer : **2**

Problem 2.1 - Linear Regression Model

Based on the plot we just made, it seems that a linear regression model could be a good modeling choice. Based on our understanding of the data from the previous subproblem, which model best describes our estimation problem?

1. $ILI = \text{intercept} + \text{coefficient} \times \text{Queries}$, where the coefficient is negative
2. $Queries = \text{intercept} + \text{coefficient} \times ILI$, where the coefficient is negative
3. $ILI = \text{intercept} + \text{coefficient} \times \text{Queries}$, where the coefficient is positive
4. $Queries = \text{intercept} + \text{coefficient} \times ILI$, where the coefficient is positive
 $\log(ILI) = \text{intercept} + \text{coefficient} \times \text{Queries}$, where the coefficient is negative

5. $\text{Queries} = \text{intercept} + \text{coefficient} \times \log(\text{ILI})$, where the coefficient is negative
6. $\log(\text{ILI}) = \text{intercept} + \text{coefficient} \times \text{Queries}$, where the coefficient is positive
7. $\text{Queries} = \text{intercept} + \text{coefficient} \times \log(\text{ILI})$, where the coefficient is positive

Answer :6

Problem 2.2 - Linear Regression Model

Let's call the regression model from the previous problem (Problem 2.1) `FluTrend1` and run it in R. Hint: to take the logarithm of a variable `Var` in a regression equation, you simply use `log(Var)` when specifying the formula to the `lm()` function.

What is the training set R-squared value for `FluTrend1` model (the "Multiple R-squared")?

```
FluTrend1 <- lm(log(ILI)~Queries, data = FluTrain)
summary(FluTrend1)
```

```
##
## Call:
## lm(formula = log(ILI) ~ Queries, data = FluTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.76003 -0.19696 -0.01657  0.18685  1.06450
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.49934    0.03041  -16.42  <2e-16 ***
## Queries      2.96129    0.09312   31.80  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2995 on 415 degrees of freedom
## Multiple R-squared:  0.709, Adjusted R-squared:  0.7083
## F-statistic: 1011 on 1 and 415 DF, p-value: < 2.2e-16
```

Answer:0.709

Problem 2.3 - Linear Regression Model

For a single variable linear regression model, there is a direct relationship between the R-squared and the correlation between the independent and the dependent variables. What is the relationship we infer from our problem? (Don't forget that you can use the `cor` function to compute the correlation between two variables.)

1. $R\text{-squared} = \text{Correlation}^2$
2. $R\text{-squared} = \log(1/\text{Correlation})$
3. $R\text{-squared} = \exp(-0.5 \times \text{Correlation})$

```
cor(log(FluTrain$ILI),FluTrain$Queries)^2
```

```
## [1] 0.7090201
```

Note that the “exp” function stands for the exponential function. The exponential can be computed in R using the function `exp()`.

Answer: 1

Problem 3.1 - Performance on the Test Set

The csv file `FluTest.csv` provides the 2012 weekly data of the ILI-related search queries and the observed weekly percentage of ILI-related physician visits. Load this data into a data frame called `FluTest`.

```
FluTest<- read.csv("data/FluTest.csv")
```

Normally, we would obtain test-set predictions from the model `FluTrend1` using the code

```
PredTest1 = predict(FluTrend1, newdata=FluTest)
```

However, the dependent variable in our model is `log(ILI)`, so `PredTest1` would contain predictions of the `log(ILI)` value. We are instead interested in obtaining predictions of the ILI value. We can convert from predictions of `log(ILI)` to predictions of ILI via exponentiation, or the `exp()` function. The new code, which predicts the ILI value, is

```
PredTest1 = exp(predict(FluTrend1, newdata=FluTest))
```

What is our estimate for the percentage of ILI-related physician visits for the week of March 11, 2012? (HINT: You can either just output `FluTest$Week` to find which element corresponds to March 11, 2012, or you can use the “which” function in R.

```
Pred <- exp(predict(FluTrend1, newdata = FluTest))
Pred[which(FluTest$Week=='2012-03-11 - 2012-03-17')]
```

```
##          11
## 2.187378
```

Problem 3.2 - Performance on the Test Set

What is the relative error between the estimate (our prediction) and the observed value for the week of March 11, 2012? Note that the relative error is calculated as

$(\text{Observed ILI} - \text{Estimated ILI}) / \text{Observed ILI}$

```
RelativeError <- (FluTest$ILI - Pred) /FluTest$ILI
RelativeError[which(FluTest$Week=='2012-03-11 - 2012-03-17')]
```

```
##          11
## 0.04623827
```

Problem 3.3 - Performance on the Test Set

What is the Root Mean Square Error (RMSE) between our estimates and the actual observations for the percentage of ILI-related physician visits, on the test set?

```
RMSE <- sqrt(mean((Pred - FluTest$ILI)^2))
RMSE
```

```
## [1] 0.7490645
```

Problem 4.1 - Training a Time Series Model

The observations in this dataset are consecutive weekly measurements of the dependent and independent variables. This sort of dataset is called a “time series.” Often, statistical models can be improved by predicting the current value of the dependent variable using the value of the dependent variable from earlier weeks. In our models, this means we will predict the ILI variable in the current week using values of the ILI variable from previous weeks.

First, we need to decide the amount of time to lag the observations. Because the ILI variable is reported with a 1- or 2-week lag, a decision maker cannot rely on the previous week’s ILI value to predict the current week’s value. Instead, the decision maker will only have data available from 2 or more weeks ago. We will build a variable called ILILag2 that contains the ILI value from 2 weeks before the current observation.

To do so, we will use the “zoo” package, which provides a number of helpful methods for time series models. While many functions are built into R, you need to add new packages to use some functions. New packages can be installed and loaded easily in R, and we will do this many times in this class. Run the following two commands to install and load the zoo package. In the first command, you will be prompted to select a CRAN mirror to use for your download. Select a mirror near you geographically.

```
lubripack::lubripack("zoo")
```

```
##
## Attaching package: 'zoo'
##
## The following objects are masked from 'package:base':
##
##   as.Date, as.Date.numeric
##
##
## Bellow Packages Successfully Installed:
##
##   zoo
## TRUE
```

Note: If ‘lubripack’ not installed, install it from [github/espanta/lubripack](https://github.com/espanta/lubripack)

After installing and loading the zoo package, run the following commands to create the ILILag2 variable in the training set:

```
ILILag2 = lag(zoo(FluTrain$ILI), -2, na.pad=TRUE)
FluTrain$ILILag2 = coredat(ILILag2)
```

In these commands, the value of -2 passed to lag means to return 2 observations before the current one; a positive value would have returned future observations. The parameter na.pad=TRUE means to add missing values for the first two weeks of our dataset, where we can't compute the data from 2 weeks earlier.

How many values are missing in the new ILILag2 variable?

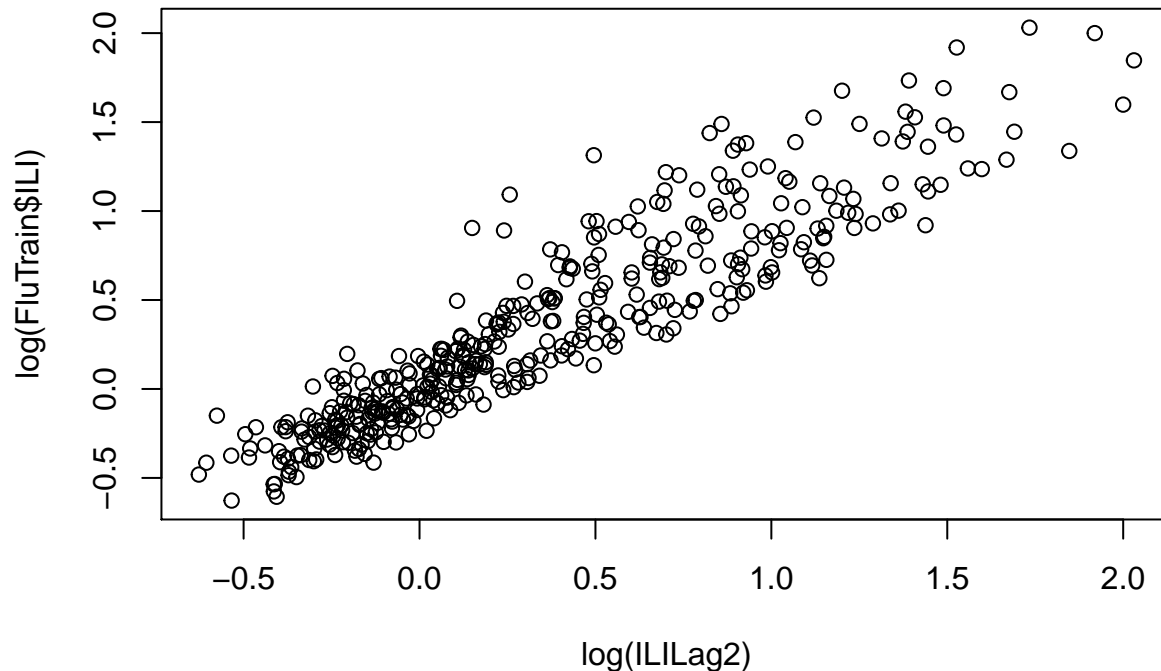
```
sum(is.na(ILILag2))
```

```
## [1] 2
```

Problem 4.2 - Training a Time Series Model

Use the plot() function to plot the log of ILILag2 against the log of ILI. Which best describes the relationship between these two variables?

```
plot(log(ILILag2), log(FluTrain$ILI))
```



1. There is a strong negative relationship between $\log(\text{ILILag2})$ and $\log(\text{ILI})$.
2. This is a weak or no relationship between $\log(\text{ILILag2})$ and $\log(\text{ILI})$.
3. There is a strong positive relationship between $\log(\text{ILILag2})$ and $\log(\text{ILI})$.

Problem 4.3 - Training a Time Series Model

Train a linear regression model on the FluTrain dataset to predict the log of the ILI variable using the Queries variable as well as the log of the ILILag2 variable. Call this model FluTrend2.

Which coefficients are significant at the $p=0.05$ level in this regression model? (Select all that apply.)

```
summary(FluTrend2 <- lm(log(ILI)~ Queries + log(ILILag2), data = FluTrain))
```

```
##
## Call:
## lm(formula = log(ILI) ~ Queries + log(ILILag2), data = FluTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.52209 -0.11082 -0.01819  0.08143  0.76785
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.24064    0.01953  -12.32  <2e-16 ***
## Queries       1.25578    0.07910   15.88  <2e-16 ***
## log(ILILag2)  0.65569    0.02251   29.14  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1703 on 412 degrees of freedom
## (2 observations deleted due to missingness)
## Multiple R-squared:  0.9063, Adjusted R-squared:  0.9059
## F-statistic: 1993 on 2 and 412 DF, p-value: < 2.2e-16
```

1. Intercept
2. Queries
3. log(ILILag2)

What is the R^2 value of the FluTrend2 model?

Answer: **0.9063**

Problem 4.4 - Training a Time Series Model

On the basis of R-squared value and significance of coefficients, which statement is the most accurate?

1. Due to overfitting, FluTrend2 is a weaker model than FluTrend1 on the training set.
2. FluTrend2 is about the same quality as FluTrend1 on the training set.
3. FluTrend2 is a stronger model than FluTrend1 on the training set.

Answer: **3**

Problem 5.1 - Evaluating the Time Series Model in the Test Set

So far, we have only added the ILILag2 variable to the FluTrain data frame. To make predictions with our FluTrend2 model, we will also need to add ILILag2 to the FluTest data frame (note that adding variables before splitting into a training and testing set can prevent this duplication of effort).

Modify the code from the previous subproblem to add an ILILag2 variable to the FluTest data frame. How many missing values are there in this new variable?

```
ILILag2 = lag(zoo(FluTest$ILI), -2, na.pad=TRUE)
FluTest$ILILag2 = coredata(ILILag2)
sum(is.na(FluTest$ILILag2))
```

```
## [1] 2
```

Problem 5.2 - Evaluating the Time Series Model in the Test Set

In this problem, the training and testing sets are split sequentially – the training set contains all observations from 2004-2011 and the testing set contains all observations from 2012. There is no time gap between the two datasets, meaning the first observation in FluTest was recorded one week after the last observation in FluTrain. From this, we can identify how to fill in the missing values for the ILILag2 variable in FluTest.

Which value should be used to fill in the ILILag2 variable for the first observation in FluTest?

1. The ILI value of the second-to-last observation in the FluTrain data frame.
2. The ILI value of the last observation in the FluTrain data frame.
3. The ILI value of the first observation in the FluTest data frame.
4. The ILI value of the second observation in the FluTest data frame.

Answer: **1**

Which value should be used to fill in the ILILag2 variable for the second observation in FluTest?

1. The ILI value of the second-to-last observation in the FluTrain data frame.
2. The ILI value of the last observation in the FluTrain data frame.
3. The ILI value of the first observation in the FluTest data frame.
4. The ILI value of the second observation in the FluTest data frame.

Answer: **2**

Problem 5.3 - Evaluating the Time Series Model in the Test Set

Fill in the missing values for ILILag2 in FluTest. In terms of syntax, you could set the value of ILILag2 in row “x” of the FluTest data frame to the value of ILI in row “y” of the FluTrain data frame with “FluTest\$ILILag2[x] = FluTrain\$ILI[y]”. Use the answer to the previous questions to determine the appropriate values of “x” and “y”. It may be helpful to check the total number of rows in FluTrain using str(FluTrain) or nrow(FluTrain).

What is the new value of the ILILag2 variable in the first row of FluTest?

```
FluTest$ILILag2[1] = FluTrain$ILI[416]  
FluTest$ILILag2[1]
```

```
## [1] 1.852736
```

What is the new value of the ILILag2 variable in the second row of FluTest?

```
FluTest$ILILag2[2] = FluTrain$ILI[417]  
FluTest$ILILag2[2]
```

```
## [1] 2.12413
```

Problem 5.4 - Evaluating the Time Series Model in the Test Set

Obtain test set predictions of the ILI variable from the FluTrend2 model, again remembering to call the `exp()` function on the result of the `predict()` function to obtain predictions for ILI instead of `log(ILI)`.

```
pred <- exp(predict(FluTrend2, newdata = FluTest))  
RMSE <- sqrt(mean((pred - FluTest$ILI)^2))  
RMSE
```

```
## [1] 0.2942029
```

What is the test-set RMSE of the FluTrend2 model?

Answer: **0.296511**

Problem 5.5 - Evaluating the Time Series Model in the Test Set

Which model obtained the best test-set RMSE?

1. FluTrend1
2. FluTrend2

Answer: **FluTrend2**

In this problem, we used a simple time series model with a single lag term. ARIMA models are a more general form of the model we built, which can include multiple lag terms as well as more complicated combinations of previous values of the dependent variable. If you're interested in learning more, check out `?arima` or the available online tutorials for these sorts of models.