

# STAT225

Oliver Baldwin Edwards, Leah Johnson, Adi Arifovic

06 May, 2020

## Predicting Hurricane Deaths

Group 2: Oliver Baldwin Edwards, Leah Johnson, Adi Arifovic

*table of contents here*

### Purpose

We hope to predict the log number of hurricane deaths based on variables such as maximum sustained windspeed, atmospheric pressure, and property damage from a dataset containing data on 94 named hurricanes that made landfall in the US mainland from 1950 through 2012. We would expect the number of deaths to increase as the maximum sustained windspeed and property damage increases, since it would make sense for more destructive hurricanes to also be more deadly.

### Data

#### Source

We plan to use the dataset “hurricNamed” from the “DAAG” R package containing data on 94 named hurricanes that made landfall in the US mainland from 1950 through 2012. It contains information on the number of deaths, the name of the hurricane, the year of the hurricane, the damage caused by the hurricane, etc. The data is sourced from multiple places and was used in a research paper claiming that hurricanes with female names did more human damage (after adjusting for the severity of the storm) than those with male names. Therefore, the data seems fairly reliable.

#### Response Variable

Our response variable is the log number of deaths that occurred due to a hurricane. The units are number of human deaths. We observe that the range of deaths are from 0 to 1836.

#### Explanatory Variables

The explanatory variables we will examine are LF.WindsMPH, LF.PressureMB, LF.times, BaseDam2014, NDAM2014, and deaths.

1. LF.WindsMPH describes the maximum sustained windspeed for each hurricane in miles per hour
2. LF.PressureMB describes the atmospheric pressure at landfall in millibars
3. LF.times describes the number of times the hurricane made landfall

4. BaseDam2014 describes the property damage caused by the hurricane in millions of 2014 US dollars
5. NDAM2014 describes the amount of damage the hurricane caused had it appeared in 2014 (no units given)
6. deaths describes the number of human deaths the hurricane caused

## Exploratory Data Analysis

### Kernel Density Estimates

The researchers first looked at kernel density estimates of our predictors, and found that `deaths`, `NDAM2014`, and `BaseDam2014` had very non-normal/skewed kernel density estimates. The researchers found that adding a `log()` transformation to `deaths`, `NDAM2014`, and `BaseDam2014` resulted in more normal kernel density estimates and minimized the effects of the outliers. (The application of a log transformation allows the researchers to fit linear regression models (OLS and JHM) to the data.) A few of the aforementioned KDEs can be seen below:

For `log_deaths`, the rectangular kernel overfits and the bandwidth selected by UCV is similarly unsuitable. The gaussian kernel looks appears to be oversmoothing. The triangular and epanechnikov kernels are similar, but the researchers decided to go with the epanechnikov kernel with bandwidth selected using BCV. The KDE looks fairly reasonable overall, and epanechnikov is generally a good choice of kernel.

### Correlation Tests

The researchers proceeded to run correlation tests to investigate the potential for positive associations between different predictors and our response variable. Here we use Kendall's  $\tau$  as it is robust against outliers. The researchers ran correlation tests with Kendall's  $\tau$  at a significance level  $\alpha = 0.05$  to test the following hypotheses: the null hypothesis of no association (independence)  $H_0 : \tau \leq 0$  and the alternative hypothesis of positive association  $H_A : \tau > 0$ . (They also tested for negative associations with the set of hypotheses  $H_0 : \tau \geq 0$  and  $H_A : \tau < 0$ .)

They found that `log_NDAM2014`, `log_BaseDam2014`, `LF.times`, and `LF.WindsMPH` were all significantly positively associated with `log_deaths`. The researchers also found that `LF.PressureMB` was significantly negatively associated with `log_deaths`. These aggregated results can be found in the table below: *(STILL NEED TO ADD THIS TABLE)*

## Model Fitting

### OLS

The researchers next proceeded to fit models to the data, beginning with ordinary least squares regression. *still need to list assumptions here for OLS?* The researchers found that the best OLS model was the model predicting `log_deaths` with simply `log_NDAM2014`. *Note that the change of `log(BaseDam2014)` has meant that none of our two predictor models are significant*

### Checking OLS Assumptions

*note we still need to plot the empirical CDFs here*

The researcher's checked the assumption of normally distributed residuals for our above selected linear regression model. To do so, the researchers performed a Kolmogorov Smirnov test using a significance level  $\alpha = 0.05$ :

**Hypotheses:**

$$H_0 : F(t) = F^*(t)$$

$$H_A : F(t) \neq F^*(t) \text{ for at least one } t$$

Where  $F(t)$  refers to the estimated CDF of the distribution of residuals of our linear model, and  $F^*(t)$  is the CDF of the normal distribution.

**Assumptions:**

1. Continuous: The test assumes that the theoretical distribution is continuous, which is reasonable.
2. Independence: This is reasonable, as there's no reason to believe the log deaths for one hurricane would affect another.
3. Parameters: This is also satisfied, we have mean = 0, and sd = 1 for the normal distribution, and neither of these is estimated.

**Test Statistic and P-value:****Decision:**

We have a  $p$ -value of  $0.7095 > 0.05$ , so we fail to reject the null hypothesis at our significance level.

**Conclusion:**

We conclude that there is insufficient evidence to suggest that the estimated CDF of the distribution of the residuals in our linear model is significantly different to the CDF of the normal distribution. Thus, our assumptions for our linear model are met.

**JHM**

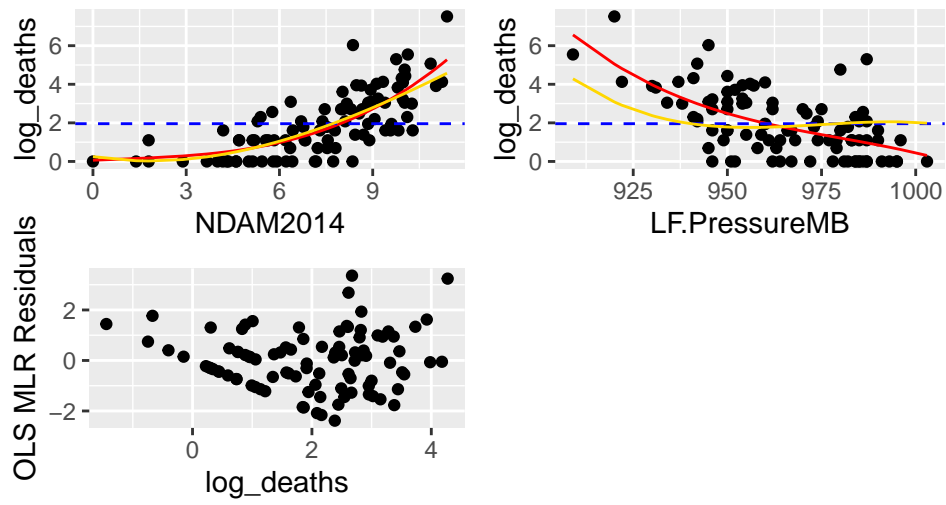
The researchers proceeded with a non-parametric approach to model fitting by using rank-based regression. They found that the best model predicting log\_deaths was again one that used only log\_NDAM2014. The researchers performed drop in dispersion tests when making this decision and found that the model containing only NDAM2014 was the best model.

**GAM**

Lastly, the researchers attempted to fit a generalized additive model (GAM) to the data. The researchers used manual forward selection to fit the best model. With an AIC of 290.8135, the researchers first found the best way to fit log\_NDAM2014 was with a b-spline with the default number of degrees of freedom.

The researchers next found that the best way to fit LF.PressureMB was with an s-spline with 5 degrees of freedom (while keeping the previous b-spline for log\_NDAM2014).

All other predictors increased the AIC, so the researchers decided on the GAM using a b-spline with the default number of degrees of freedom for log\_NDAM2014, and an s-spline with 5 degrees of freedom LF.PressureMB. The researchers plotted this, but noticed that there was a poor fit for the LF.PressureMB scatterplot. The researchers noticed that the fit for LF.PressureMB appeared to be roughly quadratic, so attempted fitting the LF.PressureMB term with a polynomial. The final GAM can then be seen below: *note we still need to overlay the OLS and JHM lines from above I think*



## Results

*still need to fill this in*