

Lecture 19

Poisson Regression

STAT 8020 Statistical Methods II
October 27, 2020

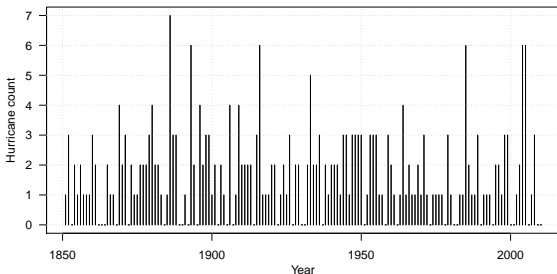
Whitney Huang
Clemson University

Daily COVID-19 Cases in South Carolina



Each day shows new cases reported since the previous day · Updated less than 19 hours ago ·
Source: [The New York Times](#) · [About this data](#)

Number of landfalling hurricanes per hurricane season



Modeling Count Data

So far we have talked about:

- Linear regression: $Y = \beta_0 + \beta_1 x + \varepsilon$, $\varepsilon \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$
- Logistic Regression: $\log\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 x$, $\pi = P(Y = 1)$

Count data

- Counts typically have a right skewed distribution
- Counts are not necessarily binary

We could use [Poisson Regression](#) to model count data

Poisson Distribution

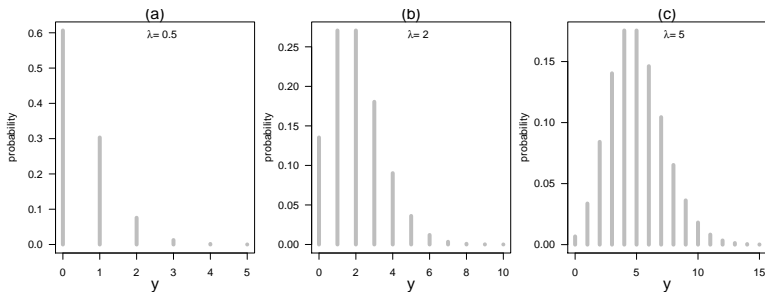
- If Y follow a Poisson distribution, then we have

$$P(Y = y) = \frac{e^{-\lambda} \lambda^y}{y!}, \quad y = 0, 1, 2, \dots,$$

where λ is the rate parameter that describe the event occurrence frequency

- $E(Y) = \text{Var}(Y) = \lambda$ if $Y \sim \text{Pois}(\lambda)$, $\lambda > 0$
- A useful model to describe the probability of a given number of events occurring in a fixed interval of time or space

Poisson Probability Mass Function



- (a), $\lambda = 0.5$: distribution gives highest probability to $y = 0$ and falls rapidly as $y \uparrow$
- (b), $\lambda = 2$: a skew distribution with longer tail on the right
- (c), $\lambda = 5$: distribution become more normally shaped

Flying-Bomb Hits on London During World War II [Clarke, 1946; Feller, 1950]

The City of London was divided into 576 small areas of one-quarter square kilometers each, and the number of areas hit exactly k times was counted. There were a total of 537 hits, so the average number of hits per area was $\frac{537}{576} = 0.9323$. The observed frequencies in the table below are remarkably close to a [Poisson distribution](#) with rate $\lambda = 0.9323$

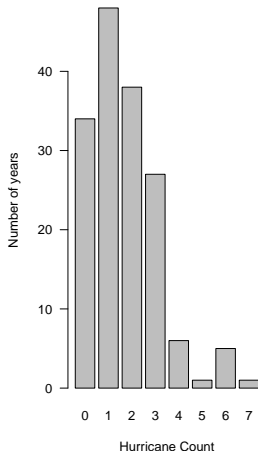
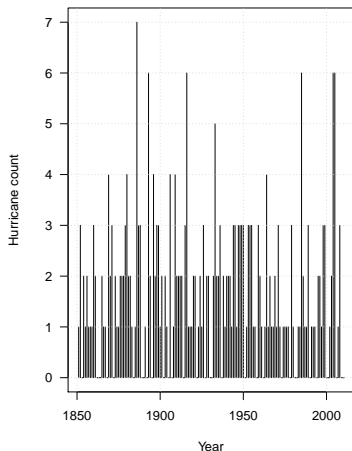
Hits	0	1	2	3	4	5+
Observed	229	211	93	35	7	1
Expected	226.7	211.4	98.5	30.6	7.1	1.6

US Hurricane Landfall points



Source: <https://www.kaggle.com/gi0vanni/analysis-on-us-hurricane-landfalls>

Number of US Landfalling Hurricanes Per Hurricane Season

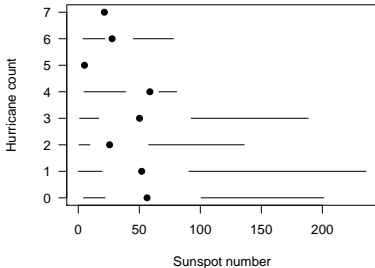
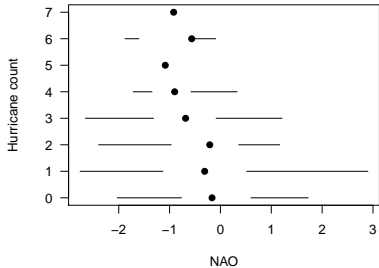
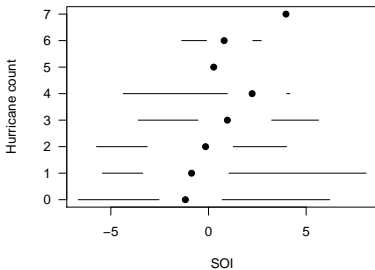
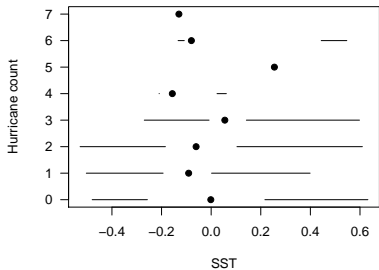


Research question: Can the variation of the annual counts be explained by some environmental variable, e.g., Southern Oscillation Index (SOI)?

Some Potentially Relevant Predictors

- Southern Oscillation Index (SOI): an indicator of wind shear
- Sea Surface Temperature (SST): an indicator of oceanic heat content
- North Atlantic Oscillation (NAO): an indicator of steering flow
- Sunspot Number (SSN): an indicator of upper air temperature

Hurricane Count vs. Environmental Variables



$$\log(\lambda) = \beta_0 + \beta_1 x_1 + \cdots + \beta_{p-1} x_{p-1}$$

$$\Rightarrow Y \sim \text{Pois}(\lambda = \exp(\beta_0 + \beta_1 x_1 + \cdots + \beta_{p-1} x_{p-1}))$$

- Model **the logarithm of the mean response** as a linear combination of the predictors
- Parameter estimation is carry out using **maximum likelihood method**
- Interpretation of β' s: every one unit increase in x_j , given that the other predictors are held constant, the λ **increases by a factor of $\exp(\beta_j)$**

Poisson Regression Model:

$$\log(\lambda_{\text{Count}}) \sim \text{SOI} + \text{NAO} + \text{SST} + \text{SSN}$$

Table: Coefficients of the Poisson regression model.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.5953	0.1033	5.76	0.0000
SOI	0.0619	0.0213	2.90	0.0037
NAO	-0.1666	0.0644	-2.59	0.0097
SST	0.2290	0.2553	0.90	0.3698
SSN	-0.0023	0.0014	-1.68	0.0928

⇒ every one unit increase in SOI, the hurricane rate increases by a factor of $\exp(0.0619) = 1.0639$ or 6.39%.

Issue with Linear Regression Fit

Linear Regression Model:

$$E(\text{Count}) \sim \text{SOI} + \text{NAO} + \text{SST} + \text{SSN}$$

Table: Coefficients of the linear regression model.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.8869	0.1876	10.06	0.0000
SOI	0.1139	0.0402	2.83	0.0053
NAO	-0.2929	0.1173	-2.50	0.0137
SST	0.4314	0.4930	0.88	0.3830
SSN	-0.0039	0.0024	-1.66	0.1000

If we use this fitted model to predict the mean hurricane count,
say SOI = -3, NAO=3, SST = 0, SSN=250

```
> predict(lmFull, newdata = data.frame(SOI = -3, NAO = 3, SST = 0, SSN = 250))
      1
-0.318065
```

This number does not make sense

Model Selection

```
> step(PoiFull)
```

```
Start: AIC=479.64
```

```
All ~ SOI + NAO + SST + SSN
```

	Df	Deviance	AIC
- SST	1	175.61	478.44
<none>		174.81	479.64
- SSN	1	177.75	480.59
- NAO	1	181.58	484.41
- SOI	1	183.19	486.02

```
Step: AIC=478.44
```

```
All ~ SOI + NAO + SSN
```

	Df	Deviance	AIC
<none>		175.61	478.44
- SSN	1	178.29	479.12
- NAO	1	183.57	484.41
- SOI	1	183.91	484.74

```
Call: glm(formula = All ~ SOI + NAO + SSN, family = "poisson", data = df)
```

```
Coefficients:
```

(Intercept)	SOI	NAO	SSN
0.584957	0.061533	-0.177439	-0.002201

```
Degrees of Freedom: 144 Total (i.e. Null); 141 Residual
```

```
Null Deviance: 197.9
```

```
Residual Deviance: 175.6            AIC: 478.4
```