

## Time Series Analysis, Regression and Forecasting

With tutorials in Python

# The Poisson Regression Model

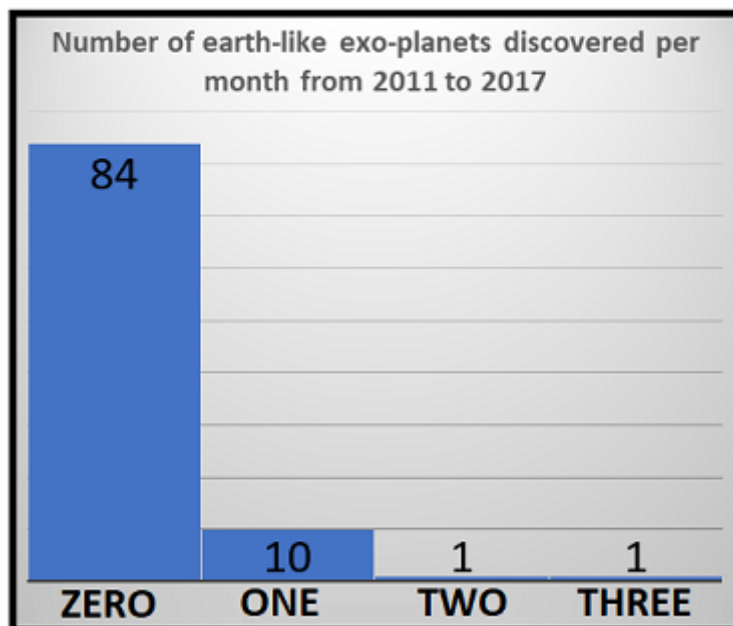
And a tutorial on Poisson regression using Python

In this section, we'll cover the following topics:

1. **Features of count based data:** Count data sets are some of the most commonly occurring data in the world. We'll look at what makes count based data different.
2. **Regression models for forecasting counts:** We'll look at **Poisson regression model** in detail. The **Negative Binomial (NB) regression model** is another commonly used model for count based data. I'll cover that in a later section.
3. **Python tutorial on Poisson regression:** I will take you through a step-by-step tutorial on how to create a Poisson regression model in Python using the **GLM class** of **statsmodels**, and how to train it on a real world data set.

## What is count based data?

Count based data contains events that occur at a certain rate. The rate of occurrence may change over time or from one observation to next. Here are some examples of count based data:



Data source: Wikipedia: [potentially habitable exoplanets](https://en.wikipedia.org/wiki/List_of_potentially_habitable_exoplanets)  
 ([https://en.wikipedia.org/wiki/List\\_of\\_potentially\\_habitable\\_exoplanets](https://en.wikipedia.org/wiki/List_of_potentially_habitable_exoplanets)  
 ). (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>))

- Number of vehicles crossing an intersection per hour,
- Number of people visiting a doctor's office per month,
- Number of earth-like planets spotted per month.

**A data set of counts has the following characteristics:**

- **Whole number data:** The data consists of non-negative integers:  $[0 \dots \infty]$  Regression techniques such as Ordinary Least Squares Regression may not be appropriate for modeling such data as OLSR works best on real numbers such as -656.0, -0.00000345, 13786.1 etc.
- **Skewed Distribution:** The data may contain a large number of data points for just a few values, thereby making the frequency distribution quite skewed. See for example above histogram.
- **Sparsity:** The data may reflect the occurrence of a rare event such as a gamma ray burst, thereby making the data sparse.
- **Rate of occurrence:** For the sake of creating a model, it can be assumed that there is a certain rate of occurrence of events  $\lambda$  that drives the generation of such data. The event rate may drift over time.

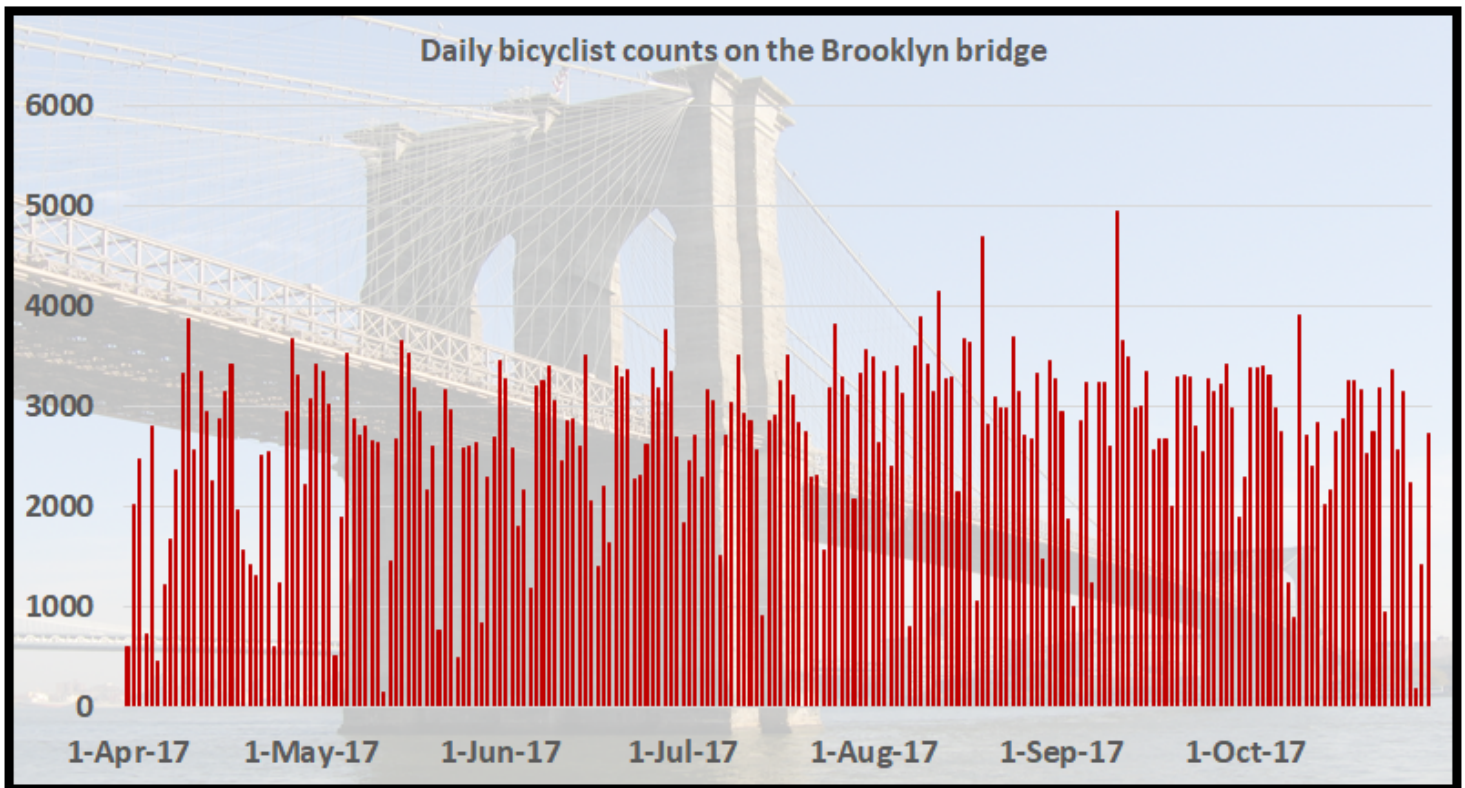
## A real world data set of counts

The following table contains counts of bicyclists traveling over various NYC bridges. The counts were measured daily from 01 April 2017 to 31 October 2017.

| Date  | Day       | High Temp (°F) | Low Temp (°F) | Precipitation | Brooklyn Bridge | Manhattan Bridge | Williamsburg Bridge | Queensboro Bridge | Total  |
|-------|-----------|----------------|---------------|---------------|-----------------|------------------|---------------------|-------------------|--------|
| 6/1   | Thursday  | 78.1           | 62.1          | 0.00          | 3,468           | 7,328            | 8,461               | 6,184             | 25,441 |
| 6/2   | Friday    | 73.9           | 60.1          | 0.01          | 3,271           | 7,007            | 7,968               | 5,293             | 23,539 |
| 6/3   | Saturday  | 72.0           | 55            | 0.01          | 2,589           | 4,510            | 6,210               | 4,084             | 17,393 |
| 6/4   | Sunday    | 68.0           | 60.1          | 0.09          | 1,805           | 3,127            | 4,023               | 3,023             | 11,978 |
| 6/5   | Monday    | 66.9           | 60.1          | 0.02          | 2,171           | 4,552            | 5,276               | 3,359             | 15,358 |
| 6/6   | Tuesday   | 55.9           | 53.1          | 0.06          | 1,193           | 3,021            | 3,807               | 2,572             | 10,593 |
| 6/7   | Wednesday | 66.9           | 54            | 0.00          | 3,211           | 7,180            | 7,632               | 5,072             | 23,095 |
| 6/8   | Thursday  | 68.0           | 59            | 0.00          | 3,253           | 7,083            | 7,778               | 5,288             | 23,402 |
| 6/9   | Friday    | 80.1           | 59            | 0.00          | 3,401           | 6,859            | 7,744               | 5,155             | 23,159 |
| 6/10  | Saturday  | 84.0           | 68            | 0.00          | 3,066           | 5,193            | 6,391               | 4,425             | 19,075 |
| 6/11  | Sunday    | 90.0           | 73            | 0.00          | 2,465           | 4,388            | 5,153               | 3,178             | 15,184 |
| 6/12  | Monday    | 91.9           | 77            | 0.00          | 2,854           | 6,265            | 7,049               | 5,032             | 21,200 |
|       |           |                |               |               |                 |                  |                     |                   |        |
| 10/26 | Thursday  | 57.0           | 53.1          | 0.00          | 2,565           | 5,217            | 5,958               | 5,011             | 18,751 |
| 10/27 | Friday    | 62.1           | 48.0          | 0.00          | 3,150           | 5,610            | 6,450               | 5,181             | 20,391 |
| 10/28 | Saturday  | 68.0           | 55.9          | 0.00          | 2,245           | 4,520            | 5,104               | 4,069             | 15,938 |
| 10/29 | Sunday    | 64.9           | 61.0          | 3.03          | 183             | 661              | 1,026               | 965               | 2,835  |
| 10/30 | Monday    | 55.0           | 46.0          | 0.25          | 1,428           | 2,966            | 3,547               | 2,924             | 10,865 |
| 10/31 | Tuesday   | 54.0           | 44.0          | 0.00          | 2,727           | 5,597            | 5,894               | 4,883             | 19,101 |

Source: Bicycle Counts for East River Bridges  
 (<https://data.cityofnewyork.us/Transportation/Bicycle-Counts-for-East-River-Bridges/gua4-p9wg>). (Data source: NYC OpenData) (Image by Author (<https://www.linkedin.com/in/sachindate/>)).

Here is a time sequenced plot of the bicyclist counts on the Brooklyn bridge:



Background image: [The Brooklyn bridge as seen from Manhattan island](https://en.wikipedia.org/wiki/Brooklyn_Bridge#/media/File:Brooklyn_Bridge_Postdlf.jpg)  
 ([https://en.wikipedia.org/wiki/Brooklyn\\_Bridge#/media/File:Brooklyn\\_Bridge\\_Postdlf.jpg](https://en.wikipedia.org/wiki/Brooklyn_Bridge#/media/File:Brooklyn_Bridge_Postdlf.jpg)).

## Regression models for counts

The **Poisson regression model** and the **Negative Binomial regression model** are two popular techniques for developing regression models for counts. Other possibilities are Ordered Logit ([https://en.wikipedia.org/wiki/Ordered\\_logit](https://en.wikipedia.org/wiki/Ordered_logit)), Ordered Probit ([https://en.wikipedia.org/wiki/Ordered\\_probit](https://en.wikipedia.org/wiki/Ordered_probit)), and Nonlinear Least Squares (<https://timeseriesreasoning.com/contents/nonlinear-least-squares-nls-regression/>) models.

## Regression strategy

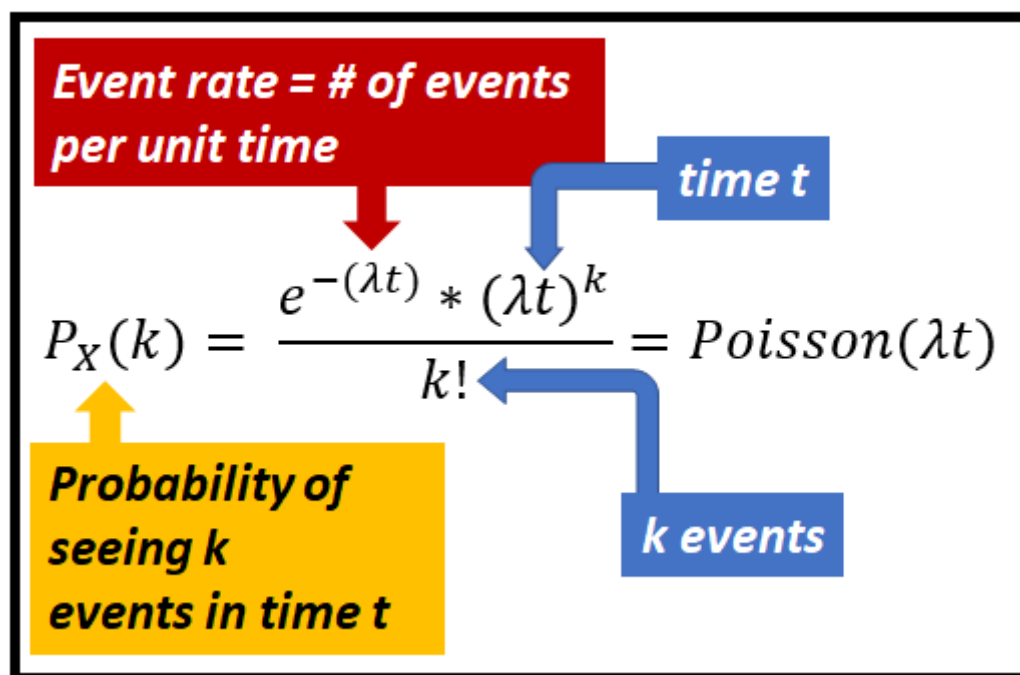
It's good practice to start with the Poisson regression model and use it as the "control" for either more complex, or less constrained models. In their book Regression Analysis of Count Data (<http://faculty.econ.ucdavis.edu/faculty/cameron/racd2/>), Cameron and Trivedi say the following:

"A sound practice is to estimate both Poisson and negative binomial models."

In this section, we'll use the Poisson regression model for regressing the bicyclist counts observed on the Brooklyn bridge, and in a following section, we'll train the Negative Binomial model (<https://timeseriesreasoning.com/contents/negative-binomial-regression-model/>) on the same data set.

## Introducing the Poisson model

The Poisson distribution has the following Probability Mass Function.



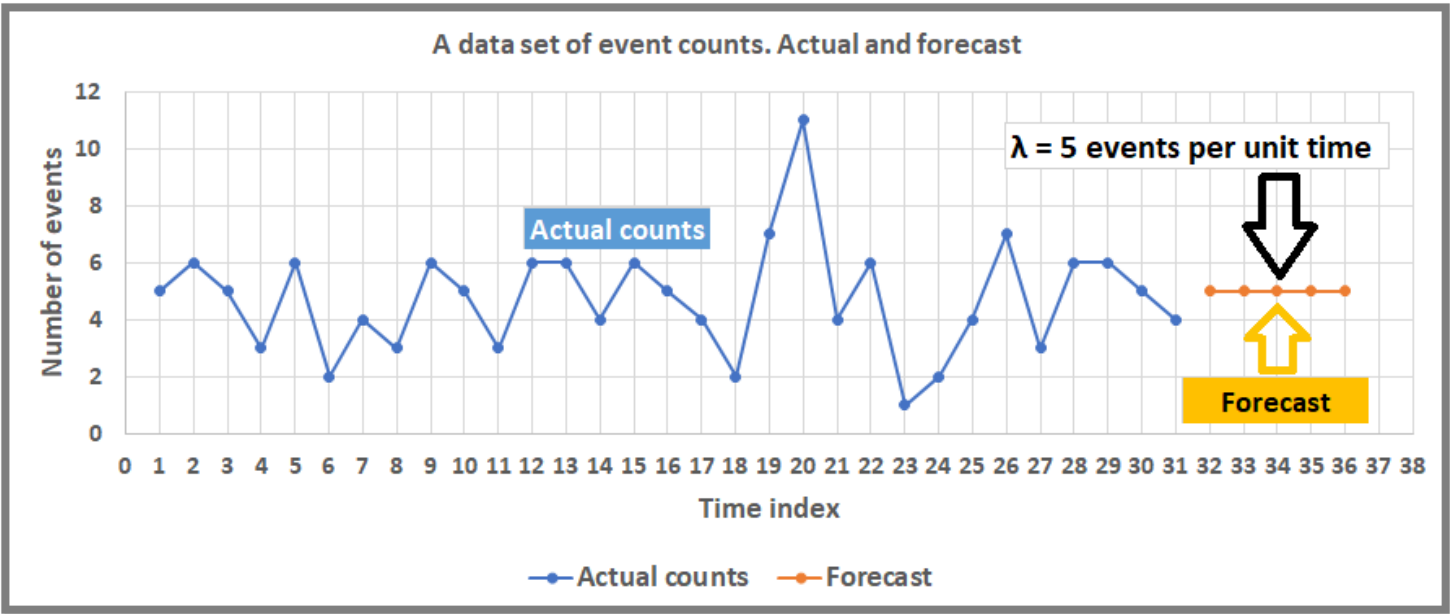
Probability of seeing  $k$  events in time  $t$ , given  
 $\lambda$  events occurring per unit time (Image  
 by [Author](#)

(<https://www.linkedin.com/in/sachindate/>).

The expected value (mean) for a Poisson distribution is  $\lambda$ . Thus in the absence of other information, one should expect to see  $\lambda$  events in any unit time interval such as 1 hour, 1 day, etc. For any interval  $t$ , one would expect to see  $\lambda t$  events.

## A Poisson regression model for *constant* $\lambda$

The following figure illustrates the constant  $\lambda$  scenario:



Actual and predicted counts for a constant rate model (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>))

The following Python code was used to generate the blue dots (actual counts in the past time steps) using a Poisson process with  $\lambda=5$ . The orange dots (predictions) are all set to the same value 5:

|    |   |
|----|---|
| 1  | import random   |
| 2  | import math   |
| 3  |   |
| 4  | _lambda = 5   |
| 5  | _num_total_arrivals = 150                               |
| 6  | _num_arrivals = 0                                       |
| 7  | _arrival_time = 0                                       |
| 8  | _num_arrivals_in_unit_time = []                         |
| 9  | _time_tick = 1  |
| 10 |   |
| 11 | print('RANDOM_N,INTER_ARRIVAL_TIME,EVENT_ARRIVAL_TIME') |
| 12 |   |
| 13 | for i in range(_num_total_arrivals):                    |
| 14 | #Get the next probability value from Uniform(0,1)       |
| 15 | p = random.random()                                     |
| 16 |   |

|    |   |
|----|---|
| 17 | # Plug it into the inverse of the CDF of Exponential(_lamnbda)  |
| 18 | _inter_arrival_time = -math.log(1.0 - p)/ _lambda   |
| 19 |   |
| 20 | # Add the inter-arrival time to the running sum   |
| 21 | _arrival_time = _arrival_time + _inter_arrival_time   |
| 22 |   |
| 23 | # Increment the number of arrival per unit time   |
| 24 | _num_arrivals = _num_arrivals + 1   |
| 25 | if _arrival_time > _time_tick:  |
| 26 | _num_arrivals_in_unit_time.append(_num_arrivals)  |
| 27 | _num_arrivals = 0   |
| 28 | _time_tick = _time_tick + 1   |
| 29 |   |
| 30 | # print it all out  |
| 31 | print(str(p)+' '+str(_inter_arrival_time)+' '+str(_arrival_time))   |
| 32 |   |
| 33 | print('\nNumber of arrivals in successive unit length intervals ==>')   |
| 34 | print(_num_arrivals_in_unit_time)   |
| 35 |   |
| 36 | print('Mean arrival rate for sample:' + str(sum(_num_arrivals_in_unit_time)/len(_num_arrivals_in_unit_time))) |

view raw poisson\_counts\_generator.py hosted with ❤ by GitHub

A Python program to generate event counts using a Poisson process

## A Poisson regression model for a **non-constant** $\lambda$

Now we get to the fun part. Let us examine a more common situation, one where  $\lambda$  can change from one observation to the next. In this case, we assume that the value of  $\lambda$  is influenced by a vector of **explanatory variables, also known as predictors, regression variables, or regressors**. We'll call this matrix of regression variables, **X**.

The job of the regression model is to fit the observed counts **y** to the matrix of regression values **X**.

In the NYC bicyclist counts data set, the regression variables are *Date*, *Day of Week*, *High Temp*, *Low Temp* and *Precipitation*. We can also introduce additional regressors such as *Month* and *Day of Month* that are derived from *Date*, and we have the liberty to drop existing regressors such as *Date*.

| Date | Day       | High Temp (°F) | Low Temp (°F) | Precipitation | Brooklyn Bridge |
|------|-----------|----------------|---------------|---------------|-----------------|
| 6/1  | Thursday  | 78.1           | 62.1          | 0.00          | 3,468           |
| 6/2  | Friday    | 73.9           | 60.1          | 0.01          | 3,271           |
| 6/3  | Saturday  | 72.0           | 55            | 0.01          | 2,589           |
| 6/4  | Sunday    | 68.0           | 60.1          | 0.09          | 1,805           |
| 6/5  | Monday    | 66.9           | 60.1          | 0.02          | 2,171           |
| 6/6  | Tuesday   | 55.9           | 53.1          | 0.06          | 1,193           |
| 6/7  | Wednesday | 66.9           | 54            | 0.00          | 3,211           |
| 6/8  | Thursday  | 68.0           | 59            | 0.00          | 3,253           |
| 6/9  | Friday    | 80.1           | 59            | 0.00          | 3,401           |
| 6/10 | Saturday  | 84.0           | 68            | 0.00          | 3,066           |
| 6/11 | Sunday    | 90.0           | 73            | 0.00          | 2,465           |
| 6/12 | Monday    | 91.9           | 77            | 0.00          | 2,854           |

(Image by [Author](#)

(<https://www.linkedin.com/in/sachindate/>).

)

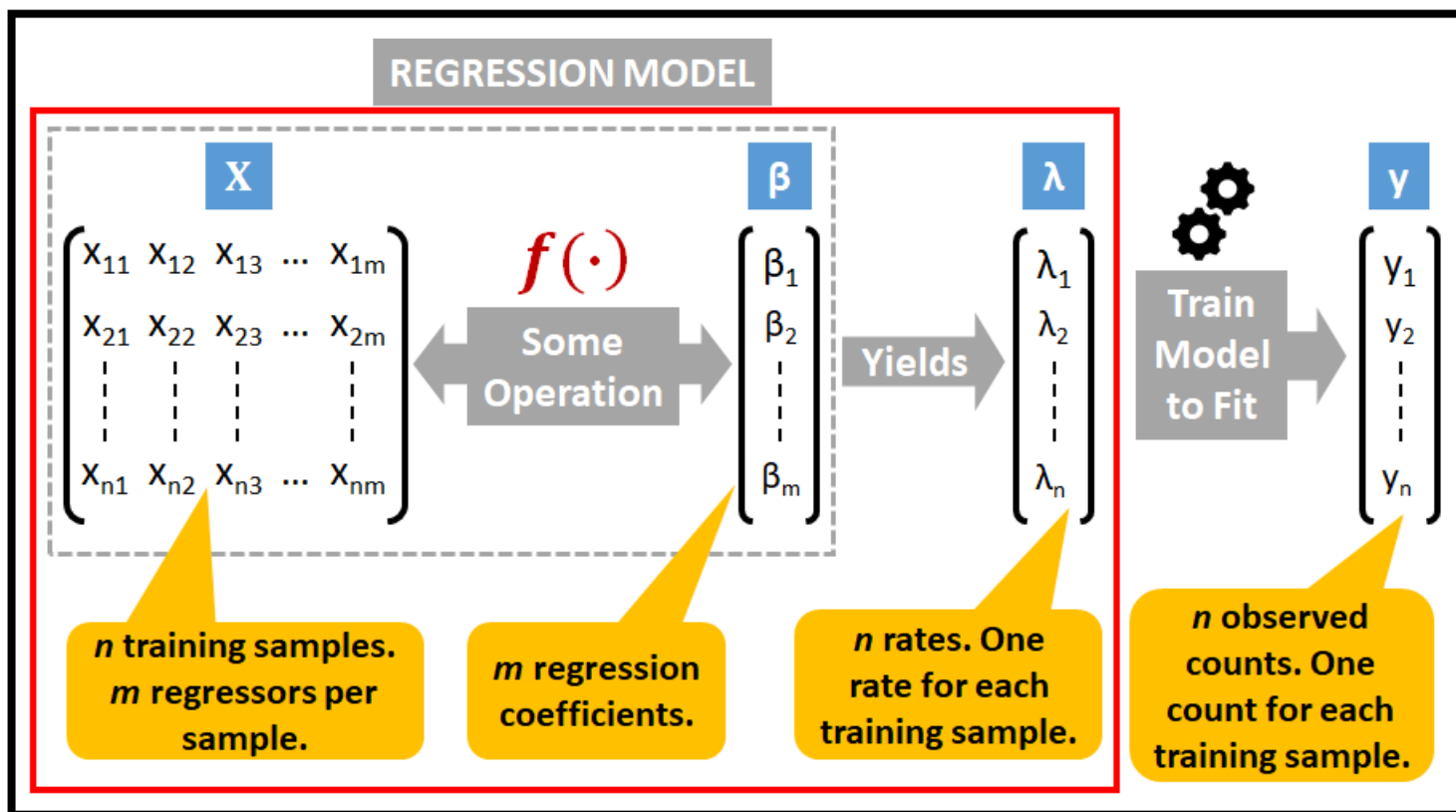
The fitting of  $\mathbf{y}$  to  $\mathbf{X}$  happens by fixing the values of a vector of regression coefficients  $\boldsymbol{\beta}$ .

In a Poisson Regression model, the event counts  $\mathbf{y}$  are assumed to be Poisson distributed, which means the probability of observing  $\mathbf{y}$  is a function of the event rate vector  $\boldsymbol{\lambda}$ .

The job of the Poisson Regression model is to fit the observed counts  $\mathbf{y}$  to the regression matrix  $\mathbf{X}$  via a link-function that expresses the rate vector  $\boldsymbol{\lambda}$  as a function of, 1) the regression coefficients  $\boldsymbol{\beta}$  and 2) the regression matrix  $\mathbf{X}$ .

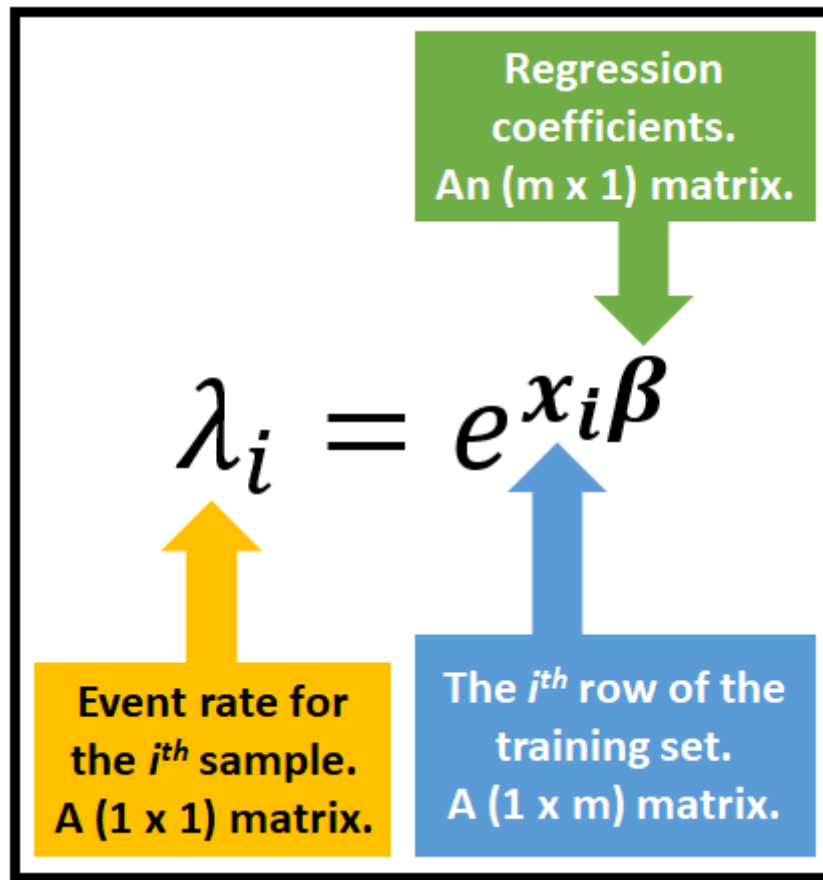
The following figure illustrates the structure of the Poisson regression model.





**Scan from left to right:** The structure of the Poisson regression model (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

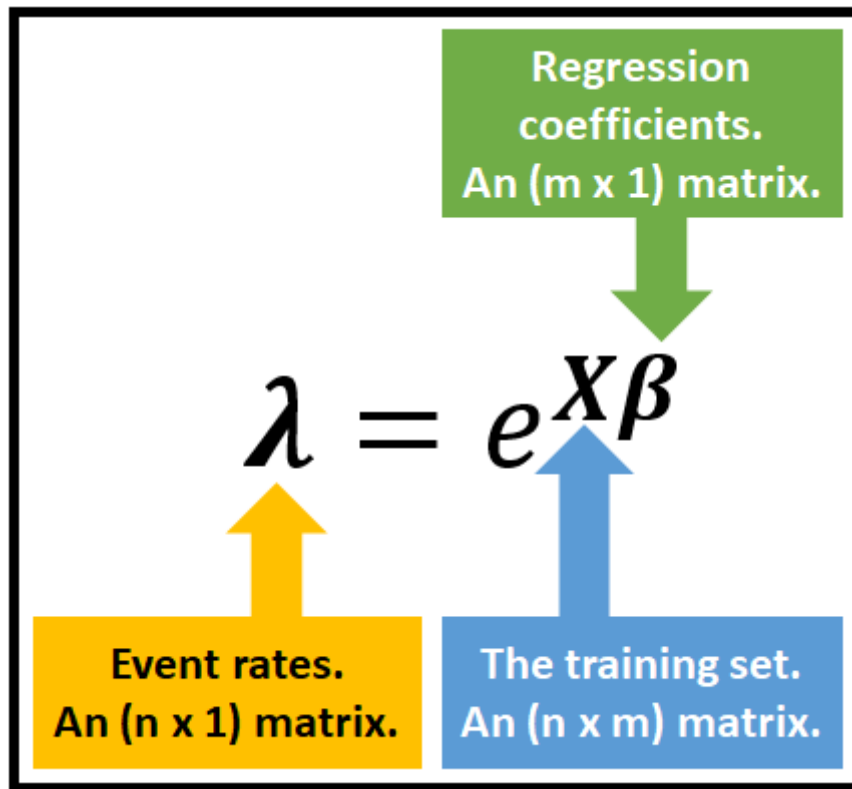
What might be a good link function  $f(\cdot)$  connecting  $\lambda$  with  $X$ ? It turns out the following exponential link-function works great:



The exponential link function of the Poisson regression model (Image by Author (<https://www.linkedin.com/in/sachindate/>)).

This link function keeps  $\lambda$  non-negative even when the regressors  $X$  or the regression coefficients  $\beta$  have negative values. This is a requirement for count based data.

In general, we have:



The exponential link function of the Poisson regression model (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

## A Formal Specification of the Poisson Regression Model

The complete specification of the Poisson regression model for count based data is given as follows:

For the  $i$ th observation in the data set denoted by  $y_i$  corresponding to the row of regression variables  $x_i$ , the probability of observing the count  $y_i$  is Poisson distributed as per the following PMF:

The diagram shows the Poisson PMF formula:  $PMF(y_i | x_i) = \frac{e^{-\lambda_i} * \lambda_i^{y_i}}{y_i!}$ . A red box at the top points to  $\lambda_i$  with the text "Event rate for the  $i^{th}$  sample". A yellow box at the bottom points to the entire formula with the text "Probability of seeing count  $y_i$  given the regression vector  $x_i$ ".

Probability of observing count  $y_i$  given  $x_i$   
 (as per the Poisson PMF formula) (Image  
 by Author  
 (<https://www.linkedin.com/in/sachindate/>).  
 )

Where the mean rate  $\lambda_i$  for the  $i^{th}$  sample is given by the exponential link function shown earlier. We reproduce it here:

The diagram shows the exponential link function:  $\lambda_i = e^{x_i \beta}$ . A blue box at the top left points to  $x_i$  with the text "Regressors for the  $i^{th}$  sample". A blue box at the top right points to  $\beta$  with the text "Regression coefficients vector". A red box at the bottom points to  $\lambda_i$  with the text "Event rate for the  $i^{th}$  sample".

The exponential link function of the Poisson  
 regression model (Image by Author  
 (<https://www.linkedin.com/in/sachindate/>).  
 )

Once the model is fully trained on the data set, the regression coefficients  $\beta$  are known, and the model is ready to make predictions. To predict the event count  $y_p$  corresponding to an input row of regressors  $x_p$  that one has observed, one uses this formula:

$$y_p = \lambda_p = e^{x_p \beta}$$

The prediction equation for the Poisson regression model (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

All of this hinges on our ability to train the model successfully so that the regression coefficients vector  $\beta$  is known.

Let's look at how this training takes place.

## Training the Poisson regression model

Training a Poisson regression model involves finding the values of the regression coefficients  $\beta$  which would make the vector of observed counts  $y$  most likely.

The technique for identifying the coefficients  $\beta$  is called **Maximum Likelihood Estimation (MLE)**.

Let's get acquainted with the technique of **MLE**.

## Understanding Maximum Likelihood Estimation (MLE)

I'll illustrate the **MLE** technique using the bicyclist counts data set. Take a look at the first few rows of this data set:

| Date | Day       | High Temp (°F) | Low Temp (°F) | Precipitation | Brooklyn Bridge |
|------|-----------|----------------|---------------|---------------|-----------------|
| 6/1  | Thursday  | 78.1           | 62.1          | 0.00          | 3,468           |
| 6/2  | Friday    | 73.9           | 60.1          | 0.01          | 3,271           |
| 6/3  | Saturday  | 72.0           | 55            | 0.01          | 2,589           |
| 6/4  | Sunday    | 68.0           | 60.1          | 0.09          | 1,805           |
| 6/5  | Monday    | 66.9           | 60.1          | 0.02          | 2,171           |
| 6/6  | Tuesday   | 55.9           | 53.1          | 0.06          | 1,193           |
| 6/7  | Wednesday | 66.0           | 54            | 0.00          | 2,000           |

Daily counts of bicyclists on the Brooklyn bridge (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

Our assumption is that the bicyclist counts shown in the red box arise from a Poisson process. Hence we can say that their probabilities of occurrence is given by the Poisson PMF. Here are the probabilities for the first 4 occurrences:

$$\begin{aligned}
 P(3468|x_1) &= \frac{e^{-\lambda_1} * \lambda_1^{3468}}{3468!} \\
 P(3271|x_2) &= \frac{e^{-\lambda_2} * \lambda_2^{3271}}{3271!} \\
 P(2589|x_3) &= \frac{e^{-\lambda_3} * \lambda_3^{2589}}{2589!} \\
 P(1805|x_4) &= \frac{e^{-\lambda_4} * \lambda_4^{1805}}{1805!}
 \end{aligned}$$

Probabilities of observing the bicyclist counts for the first few occurrences given corresponding regression vectors (Image by [Author](https://timeseriesreasoning.com/contents/poisson-regression-model/)

(<https://www.linkedin.com/in/sachindate/>)

We can similarly calculate the probabilities for all  $n$  counts observed in the training set.

Note that in the above formulae,  $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$  are calculated using the link function as follows:

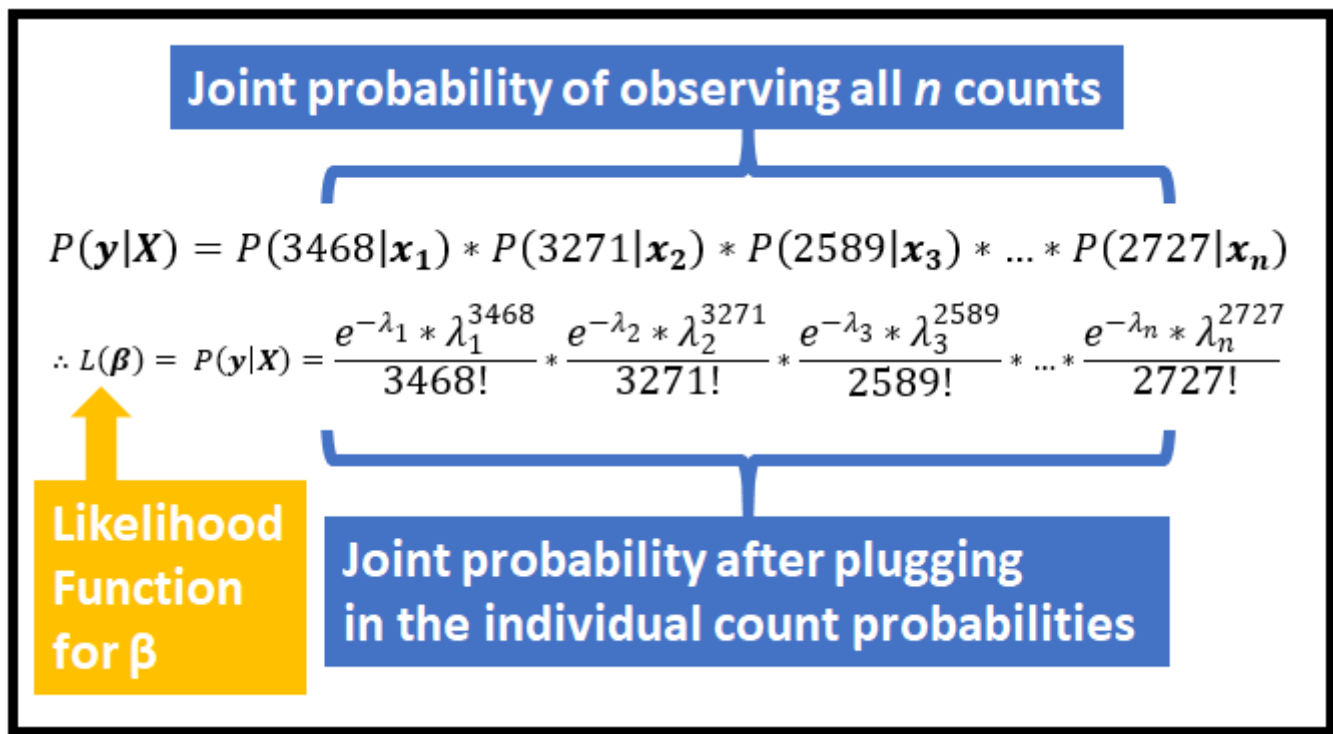
$$\begin{aligned}\lambda_1 &= e^{x_1\beta} \\ \lambda_2 &= e^{x_2\beta} \\ \lambda_3 &= e^{x_3\beta} \\ \lambda_4 &= e^{x_4\beta}\end{aligned}$$

The event rates corresponding to the counts  
for the first few days (Image by [Author](#)  
(<https://www.linkedin.com/in/sachindate/>)

Where  $x_1, x_2, x_3, x_4$  are the first 4 rows of the regression matrix.

The probability of occurrence of the entire set of  $n$  counts  $y_1, y_2, \dots, y_n$  in the training set is the joint probability of occurrence of the individual counts.

The counts  $y$  are Poisson distributed,  $y_1, y_2, \dots, y_n$  are independent random variables, given correspondingly  $x_1, x_2, \dots, x_n$ . Hence the joint probability of occurrence of  $y_1, y_2, \dots, y_n$  can be expressed as a simple multiplication of the individual probabilities. Here is how the joint probability looks like for the entire training set:



The **likelihood function**  $L(\boldsymbol{\beta})$  expressed as a joint probability mass function (Image by Author (<https://www.linkedin.com/in/sachindate/>)).

Let's recollect that  $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$  are linked to the regression vectors  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n$  via the regression coefficients  $\boldsymbol{\beta}$ .

What value of  $\boldsymbol{\beta}$  will make the given set of observed counts  $\mathbf{y}$  most likely? It is the value of  $\boldsymbol{\beta}$  for which the joint probability shown in the above equation achieves the maximum value. In other words, it is the value of  $\boldsymbol{\beta}$  for which the rate of change of the joint probability function w.r.t.  $\boldsymbol{\beta}$  is 0. Put another way, it is the solution of the equation obtained from differentiating the joint probability equation w.r.t.  $\boldsymbol{\beta}$  and setting this differential equation to 0.

It is easier to differentiate the **logarithm** of the joint probability equation than the original equation. The solution to the logged equation yields the same optimal value of  $\boldsymbol{\beta}$ .

This logarithmic equation is called the **log-likelihood function**. For the Poisson regression, the log-likelihood function is given by the following equation:



$$\ln L(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i \mathbf{x}_i \boldsymbol{\beta} - e^{\mathbf{x}_i \boldsymbol{\beta}} - \ln y_i!)$$

log-likelihood function for the Poisson regression model (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

The above equation is obtained by taking the natural logarithm of both sides of the joint probability function shown earlier, after substituting the  $\lambda_i$  with  $\exp(\mathbf{x}_i \boldsymbol{\beta})$ .

As mentioned earlier, we differentiate this log-likelihood equation w.r.t.  $\boldsymbol{\beta}$ , and set it to zero. This operation gives us the following equation:

$$\sum_{i=1}^n (y_i - e^{\mathbf{x}_i \boldsymbol{\beta}}) \mathbf{x}_i = \mathbf{0}$$

The Poisson MLE for  $\boldsymbol{\beta}$  is the solution to this equation (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

Solving this equation for the regression coefficients  $\boldsymbol{\beta}$  will yield the **Maximum Likelihood Estimate (MLE)** for  $\boldsymbol{\beta}$ .

To solve the above equation one uses an iterative method such as Iteratively Reweighted Least Squares (IRLS) ([https://en.wikipedia.org/wiki/Iteratively\\_reweighted\\_least\\_squares](https://en.wikipedia.org/wiki/Iteratively_reweighted_least_squares)). In practice, one does not solve this equation by hand. Instead, you use statistical software such as the Python **statsmodels** —

package which will do all the calculations for you while training the Poisson regression model on your data set.

## Summary of steps for performing Poisson Regression

In summary, here are the steps for performing a Poisson Regression on a count based data set:

1. First, make sure that your data set contains counts. One way to tell is that it contains only non-negative integer values that represent the number of occurrences of some event during some interval. In the bicyclist counts data set, it is the count of bicyclists crossing the Brooklyn bridge per day.
2. Find out (or guess) the regression variables that will influence the observed counts. In the bicyclist counts data set the regression variables are *Day of Week*, *Min Temp*, *Max Temp*, *Precipitation* etc.
3. Carve out a training data set that your regression model will train on, and a test data set that should keep aside. Do *not* train the model on the test data.
4. Use a suitable statistical software such as the **Python statsmodels package** to configure and fit the Poisson Regression model on the training data set.
5. Test the performance of the model by running it on the test data set so as to generate predicted counts. Compare them with the actual counts in the test data set.
6. Use a goodness-of-fit measure to determine how well your model has trained on the training data set.

## How to train a Poisson Regression Model in Python

Let's put into practice what we have learnt. The Python **statsmodels** (<https://www.statsmodels.org/stable/index.html>) package has excellent support for doing Poisson regression.

Let's use the Brooklyn bridge bicyclist counts data set. You can pick up the data set from [here](https://gist.github.com/sachinsdate/c17931a3f000492c1c42cf78bf4ce9fe) (<https://gist.github.com/sachinsdate/c17931a3f000492c1c42cf78bf4ce9fe>).

Our goal is to build a Poisson regression model for the observed bicyclist counts  $y$ . We will use the trained model to predict daily counts of bicyclists on the Brooklyn bridge that the model has not seen during training.

| Regression variables matrix <b>X</b> |           |                |               |               | Observed counts vector <b>y</b> |
|--------------------------------------|-----------|----------------|---------------|---------------|---------------------------------|
| Date                                 | Day       | High Temp (°F) | Low Temp (°F) | Precipitation | Brooklyn Bridge                 |
| 6/1                                  | Thursday  | 78.1           | 62.1          | 0.00          | 3,468                           |
| 6/2                                  | Friday    | 73.9           | 60.1          | 0.01          | 3,271                           |
| 6/3                                  | Saturday  | 72.0           | 55            | 0.01          | 2,589                           |
| 6/4                                  | Sunday    | 68.0           | 60.1          | 0.09          | 1,805                           |
| 6/5                                  | Monday    | 66.9           | 60.1          | 0.02          | 2,171                           |
| 6/6                                  | Tuesday   | 55.9           | 53.1          | 0.06          | 1,193                           |
| 6/7                                  | Wednesday | 66.9           | 54            | 0.00          | 3,211                           |
| 6/8                                  | Thursday  | 68.0           | 59            | 0.00          | 3,253                           |
| 6/9                                  | Friday    | 80.1           | 59            | 0.00          | 3,401                           |
| 6/10                                 | Saturday  | 84.0           | 68            | 0.00          | 3,066                           |
| 6/11                                 | Sunday    | 90.0           | 73            | 0.00          | 2,465                           |
| 6/12                                 | Monday    | 91.9           | 77            | 0.00          | 2,854                           |
|                                      |           |                |               |               |                                 |
| 10/26                                | Thursday  | 57.0           | 53.1          | 0.00          | 2,565                           |
| 10/27                                | Friday    | 62.1           | 48.0          | 0.00          | 3,150                           |
| 10/28                                | Saturday  | 68.0           | 55.9          | 0.00          | 2,245                           |
| 10/29                                | Sunday    | 64.9           | 61.0          | 3.03          | 183                             |
| 10/30                                | Monday    | 55.0           | 46.0          | 0.25          | 1,428                           |
| 10/31                                | Tuesday   | 54.0           | 44.0          | 0.00          | 2,727                           |

Daily count of bicyclists on the Brooklyn bridge (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>))

We'll start by importing all the required packages.

```
import pandas as pd
from patsy import dmatrices
import numpy as np
import statsmodels.api as sm
import matplotlib.pyplot as plt
```

Create a pandas DataFrame for the counts data set.

```
df = pd.read_csv('nyc_bb_bicyclist_counts.csv', header=0,
infer_datetime_format=True, parse_dates=[0], index_col=[0])
```

We'll add a few derived regression variables to the X matrix.

```
ds = df.index.to_series()
df['MONTH'] = ds.dt.month
df['DAY_OF_WEEK'] = ds.dt.dayofweek
df['DAY'] = ds.dt.day
```

We will not use the *Date* variable as a regressor since it contains an absolute date value but we don't need to do anything special to drop *Date* as it is already consumed as the index of the pandas DataFrame. So it will not be available to us in the X matrix.

Let's create the training and testing data sets.

```
mask = np.random.rand(len(df)) < 0.8
df_train = df[mask]
df_test = df[~mask]
print('Training data set length='+str(len(df_train)))
print('Testing data set length='+str(len(df_test)))
```

Setup the regression expression in patsy notation. We are telling patsy that BB\_COUNT is our dependent variable and it depends on the regression variables: DAY, DAY\_OF\_WEEK, MONTH, HIGH\_T, LOW\_T and PRECIP.

```
expr = """BB_COUNT ~ DAY + DAY_OF_WEEK + MONTH + HIGH_T + LOW_T +
PRECIP"""
```

Set up the X and y matrices for the training and testing data sets. patsy makes this really simple.

```
y_train, X_train = dmatrices(expr, df_train, return_type='dataframe')
y_test, X_test = dmatrices(expr, df_test, return_type='dataframe')
```

Using the statsmodels GLM class, train the Poisson regression model on the training data set.

```
poisson_training_results = sm.GLM(y_train, X_train,
family=sm.families.Poisson()).fit()
```

Print the training summary.

```
print(poisson_training_results.summary())
```

This prints out the following:

Generalized Linear Model Regression Results

Dep. Variable: BB\_COUNT

No. Observations: 170

Model: GLM

Df Residuals: 163

Model Family: Poisson

Df Model: 6

Link Function: log

Scale: 1.0000

Method: IRLS

Log-Likelihood: -12335.

Date: Fri, 20 Sep 2019

Deviance: 23030.

Time: 20:11:31

Pearson chi2: 2.33e+04

No. Iterations: 5

Covariance Type: nonrobust

|             | coef      | std err | z       | P> z  | [0.025 | 0.975] |
|-------------|-----------|---------|---------|-------|--------|--------|
| Intercept   | 6.9515    | 0.012   | 564.893 | 0.000 | 6.927  | 6.976  |
| DAY         | 9.858e-06 | 0.000   | 0.057   | 0.955 | -0.000 | 0.000  |
| DAY_OF_WEEK | -0.0182   | 0.001   | -24.531 | 0.000 | -0.020 | -0.017 |
| MONTH       | 0.0176    | 0.001   | 22.469  | 0.000 | 0.016  | 0.019  |
| HIGH_T      | 0.0246    | 0.000   | 73.508  | 0.000 | 0.024  | 0.025  |
| LOW_T       | -0.0147   | 0.000   | -39.611 | 0.000 | -0.015 | -0.014 |
| PRECIP      | -0.7578   | 0.008   | -93.439 | 0.000 | -0.774 | -0.742 |

Regression Coefficients  $\beta$

p value

95% confidence interval

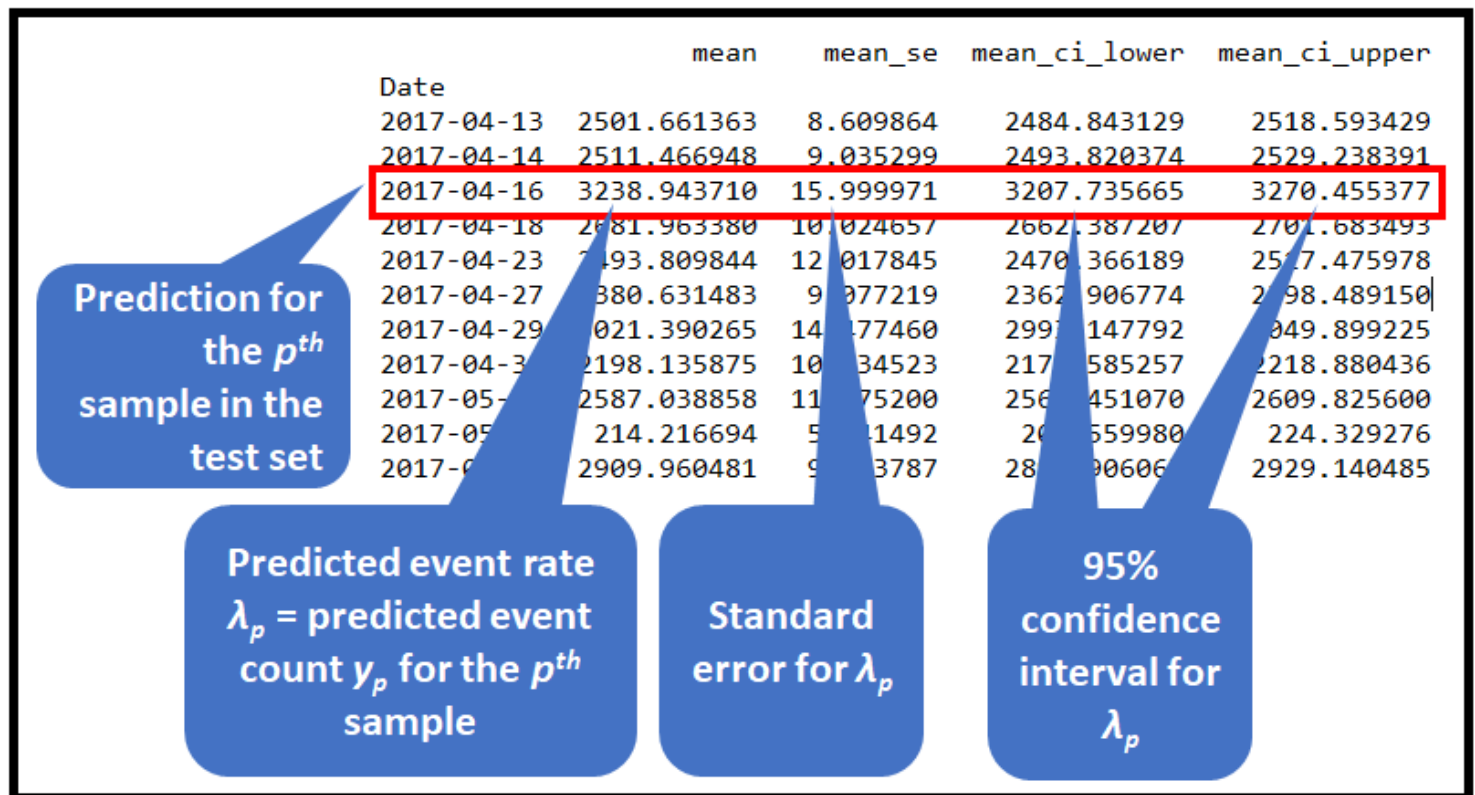
Training summary for the Poisson regression model (Image by [Author](#)

(<https://www.linkedin.com/in/sachindate/>)

So how well did our model do? Let's make some **predictions** on the test data set.

```
poisson_predictions = poisson_training_results.get_prediction(X_test)
#summary_frame() returns a pandas DataFrame
predictions_summary_frame = poisson_predictions.summary_frame()
print(predictions_summary_frame)
```

Here are the first few rows of the output:



First few rows of output from  
`poisson_predictions.summary_frame()`

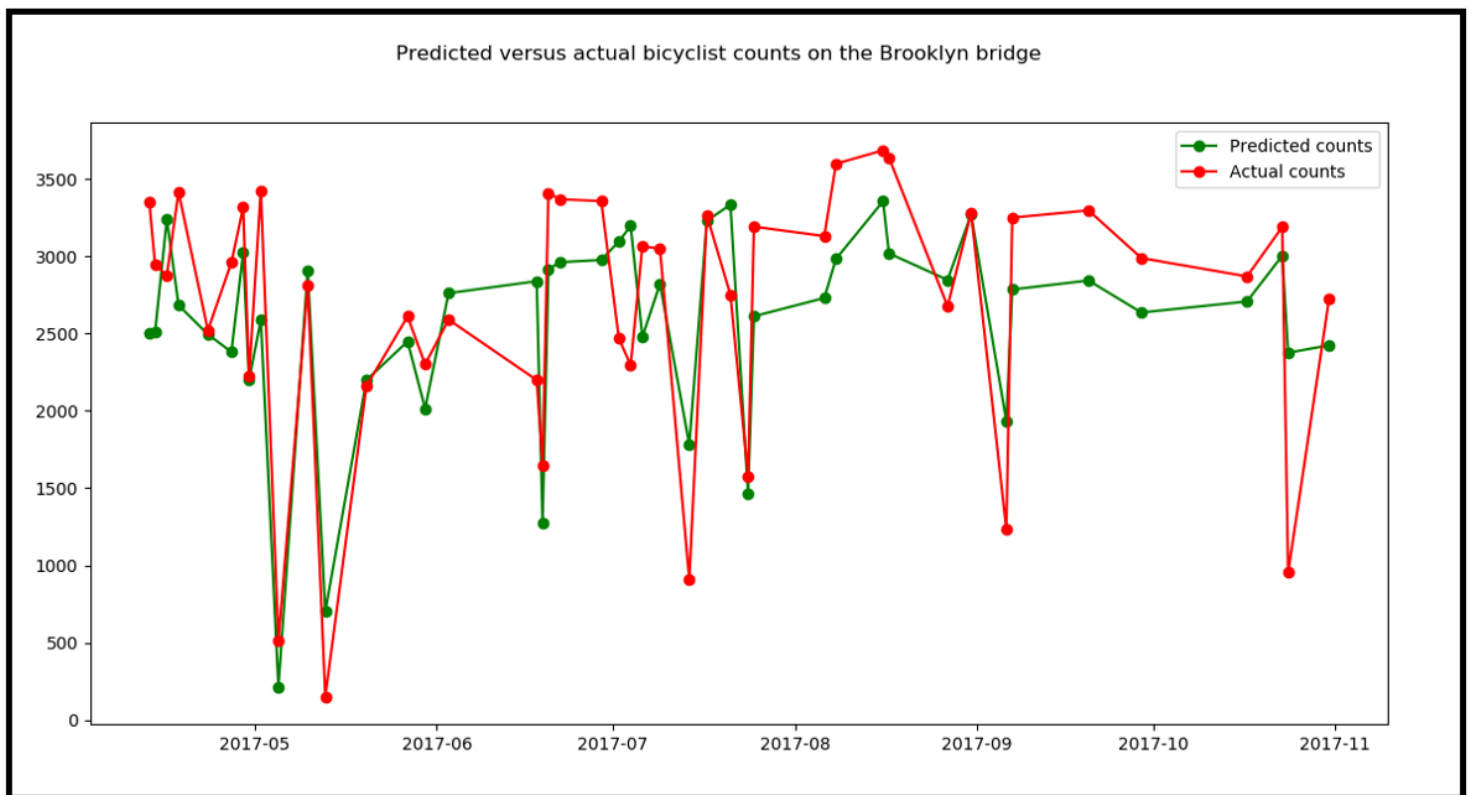
(Image by [Author](#))

(<https://www.linkedin.com/in/sachindate/>)

Let's plot the predicted counts versus the actual counts for the test data.

```
predicted_counts=predictions_summary_frame['mean']
actual_counts = y_test['BB_COUNT']
fig = plt.figure()
fig.suptitle('Predicted versus actual bicyclist counts on the Brooklyn bridge')
predicted, = plt.plot(X_test.index, predicted_counts, 'go-',
label='Predicted counts')
actual, = plt.plot(X_test.index, actual_counts, 'ro-', label='Actual counts')
plt.legend(handles=[predicted, actual])
plt.show()
```

Here's the output:



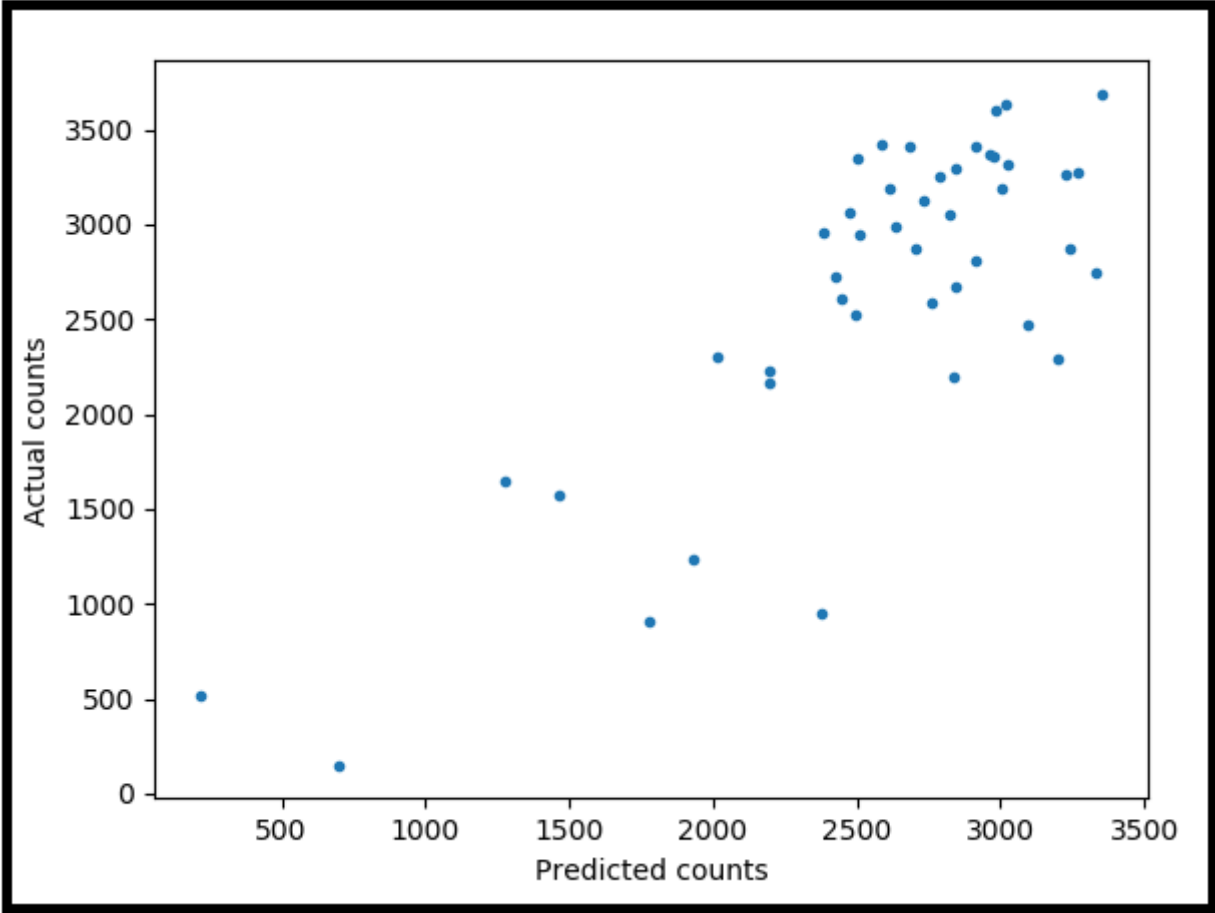
Predicted versus actual bicyclist counts on the Brooklyn bridge (Image by [Author](https://www.linkedin.com/in/sachindate/) (<https://www.linkedin.com/in/sachindate/>)).

The model seems to be more or less tracking the trend in the actual counts although **in many cases its predictions are way off the actual value.**

Let's also plot Actual versus Predicted counts.

```
plt.clf()
fig = plt.figure()
fig.suptitle('Scatter plot of Actual versus Predicted counts')
plt.scatter(x=predicted_counts, y=actual_counts, marker='.')
plt.xlabel('Predicted counts')
plt.ylabel('Actual counts')
plt.show()
```

Here is the plot:



Scatter plot of Actual versus Predicted counts  
(Image by [Author](#)  
(<https://www.linkedin.com/in/sachindate/>).  
)

Here is the complete source code for doing Poisson regression using Python:

|    |  |
|----|--|
| 1  | import pandas as pd  |
| 2  | from patsy import dmatrices  |
| 3  | import numpy as np   |
| 4  | import statsmodels.api as sm   |
| 5  | import matplotlib.pyplot as plt  |
| 6  |  |
| 7  |  |
| 8  | # Create a pandas DataFrame for the counts data set.   |
| 9  | df = pd.read_csv('nyc_bb_bicyclist_counts.csv', header=0, infer_datetime_format=True, parse_date |
| 10 |  |
| 11 | # Add a few derived regression variables.  |
| 12 | ds = df.index.to_series()  |
| 13 | df['MONTH'] = ds.dt.month  |
| 14 | df['DAY_OF_WEEK'] = ds.dt.dayofweek  |



|    |  |
|----|--|
| 15 | <code>df['DAY'] = ds.dt.day</code>   |
| 16 |  |
| 17 | <code># Create the training and testing data sets.</code>  |
| 18 | <code>mask = np.random.rand(len(df)) &lt; 0.8</code>   |
| 19 | <code>df_train = df[mask]</code>   |
| 20 | <code>df_test = df[~mask]</code>   |
| 21 | <code>print("Training data set length="+str(len(df_train)))</code>   |
| 22 | <code>print("Testing data set length="+str(len(df_test)))</code>   |
| 23 |  |
| 24 | <code># Setup the regression expression in patsy notation. We are telling patsy that BB_COUNT is our de</code> |
| 25 | <code># it depends on the regression variables: DAY, DAY_OF_WEEK, MONTH, HIGH_T, LOW_T and F</code>            |
| 26 | <code>expr = """"BB_COUNT ~ DAY + DAY_OF_WEEK + MONTH + HIGH_T + LOW_T + PRECIP""""</code>                     |
| 27 |  |
| 28 | <code># Set up the X and y matrices</code>   |
| 29 | <code>y_train, X_train = dmatrices(expr, df_train, return_type='dataframe')</code>                             |
| 30 | <code>y_test, X_test = dmatrices(expr, df_test, return_type='dataframe')</code>                                |
| 31 |  |
| 32 | <code># Using the statsmodels GLM class, train the Poisson regression model on the training data set.</code>   |
| 33 | <code>poisson_training_results = sm.GLM(y_train, X_train, family=sm.families.Poisson()).fit()</code>           |
| 34 |  |
| 35 | <code># Print the training summary.</code>   |
| 36 | <code>print(poisson_training_results.summary())</code>   |
| 37 |  |
| 38 | <code># Make some predictions on the test data set.</code>   |
| 39 | <code>poisson_predictions = poisson_training_results.get_prediction(X_test)</code>                             |
| 40 | <code># .summary_frame() returns a pandas DataFrame</code>   |
| 41 | <code>predictions_summary_frame = poisson_predictions.summary_frame()</code>                                   |
| 42 | <code>print(predictions_summary_frame)</code>  |
| 43 |  |
| 44 | <code>predicted_counts=predictions_summary_frame['mean']</code>  |
| 45 | <code>actual_counts = y_test['BB_COUNT']</code>  |
| 46 |  |
| 47 | <code># Mlot the predicted counts versus the actual counts for the test data.</code>                           |
| 48 | <code>fig = plt.figure()</code>  |
| 49 | <code>fig.suptitle('Predicted versus actual bicyclist counts on the Brooklyn bridge')</code>                   |
| 50 | <code>predicted, = plt.plot(X_test.index, predicted_counts, 'go-', label='Predicted counts')</code>            |
| 51 | <code>actual, = plt.plot(X_test.index, actual_counts, 'ro-', label='Actual counts')</code>                     |
| 52 | <code>plt.legend(handles=[predicted, actual])</code>   |
| 53 | <code>plt.show()</code>  |
| 54 |  |
| 55 | <code># Show scatter plot of Actual versus Predicted counts</code>   |
| 56 | <code>plt.clf()</code>   |

|    |   |
|----|---|
| 57 | <code>fig = plt.figure()</code>   |
| 58 | <code>fig.suptitle('Scatter plot of Actual versus Predicted counts')</code> |
| 59 | <code>plt.scatter(x=predicted_counts, y=actual_counts, marker='.')</code>   |
| 60 | <code>plt.xlabel('Predicted counts')</code>                                 |
| 61 | <code>plt.ylabel('Actual counts')</code>                                    |
| 62 | <code>plt.show()</code>   |

view raw `poisson_regression.py` hosted with ❤ by GitHub

Poisson Regression Model

## Goodness-of-fit of the Poisson regression model

Recollect that both the expected value (i.e. mean) and the variance, of the Poisson distribution is  $\lambda$ . This rather strict condition is violated by most real-world data.

A common source of failure of the Poisson regression model is that the data does not satisfy the *mean = variance* criterion imposed by the Poisson distribution.

The `summary()` method on the `statsmodels GLMResults`

([https://www.statsmodels.org/devel/generated/statsmodels.genmod.generalized\\_linear\\_model.GLMResults.html](https://www.statsmodels.org/devel/generated/statsmodels.genmod.generalized_linear_model.GLMResults.html)) class shows a couple of useful goodness-of-fit statistics to help you evaluate whether your Poisson regression model was able to successfully fit the training data. Let's look at their values:

| Generalized Linear Model Regression Results |                  |                   |           |
|---|------------------|-------------------|-----------|
| =====                                       |                  |                   |           |
| Dep. Variable:                              | BB_COUNT         | No. Observations: | 170       |
| Model:                                      | GLM              | Df Residuals:     | 163       |
| Model Family:                               | Poisson          | Df Model:         | 6         |
| Link Function:                              | log              | Scale:            | 1.0000    |
| Method:                                     | IRLS             | Log-Likelihood:   | -12335.   |
| Date:                                       | Fri, 20 Sep 2019 | Deviance:         | 23030.    |
| Time:                                       | 20:11:31         | Pearson chi2:     | 2.33e+04  |
| No. Iterations:                             | 5                | Covariance type:  | nonrobust |

Training summary for the Poisson regression model (Image by [Author](#))

(<https://www.linkedin.com/in/sachindate/>)

The reported values of Deviance and Pearson chi-squared are very large. A good fit is virtually impossible given these values. To make a quantitative determination of the goodness-of-fit at some confidence level, say 95% ( $p=0.05$ ), we look up the value in the  $\chi^2$  table for  $p=0.05$  and Degrees of freedom of residuals=163. (DF Residuals = No. Observations *minus* DF model]). We compare this Chi-Squared value with the observed statistic, in this case, the Deviance or the Pearson's chi-squared value reported in GLMResults. We find that at  $p=0.05$  and DF Residuals = 163, the chi-squared value from a standard Chi-Squared table (<https://www.medcalc.org/manual/chi-square-table.php>) is 193.791 which is much smaller than the reported statistic of 23030 and 23300. Hence as per this test, the Poisson regression model, in spite of demonstrating an 'okay' *visual* fit for the test data set, has fit the training data rather poorly.

## Conclusion and next steps

For counts based data, a useful starting point is the Poisson regression model. One can then compare its performance with other popular counts based models, such as:

1. The **The Negative Binomial regression model** (<https://timeseriesreasoning.com/contents/negative-binomial-regression-model/>) which does not make the *mean = variance* assumption about the data.
2. **Generalized Poisson regression models** (<https://timeseriesreasoning.com/contents/generalized-poisson-regression-model/>) which also work well with over-dispersed or under-dispersed data sets.
3. A **Zero Inflated Poisson model** (<https://timeseriesreasoning.com/contents/zero-inflated-poisson-regression-model/>) if you suspect that your data contains excess zeros i.e. many more zeroes than what the regular Poisson model can explain

Happy modeling!

## Related

[Getting to Know The Poisson Process And The Poisson Probability Distribution](https://timeseriesreasoning.com/contents/poisson-process/)  
(<https://timeseriesreasoning.com/contents/poisson-process/>).

## Citations and Copyrights

## Data set

Bicycle Counts for East River Bridges. (<https://data.cityofnewyork.us/Transportation/Bicycle-Counts-for-East-River-Bridges/gua4-p9wg>) Daily total of bike counts conducted monthly on the Brooklyn Bridge, Manhattan Bridge, Williamsburg Bridge, and Queensboro Bridge. From NYC Open Data under Terms of Use. **Curated data set for download.** (<https://gist.github.com/sachinsdate/c17931a3f000492c1c42cf78bf4ce9fe>)

## Book

Cameron A. C. and Trivedi P. K., Regression Analysis of Count Data (<http://faculty.econ.ucdavis.edu/faculty/cameron/racd2/>), Second Edition, Econometric Society Monograph No. 53, Cambridge University Press, Cambridge, May 2013.

## Images

All images are copyright Sachin Date (<https://www.linkedin.com/in/sachindate/>) under CC-BY-NC-SA (<https://creativecommons.org/licenses/by-nc-sa/4.0/>), unless a different source and copyright are mentioned underneath the image.

**PREVIOUS:** The Quantile Regression Model (<https://timeseriesreasoning.com/contents/introduction-to-the-quantile-regression-model/>).

**NEXT:** The Negative Binomial Regression Model (<https://timeseriesreasoning.com/contents/negative-binomial-regression-model/>).

**UP:** Table of Contents (<https://timeseriesreasoning.com/>).

