

# HW2a

Alan Wu

MKTG 7760: Probability Models in Marketing

February 11, 2026

## 1 Question 1

The HW2A dataset has a spreadsheet containing Internet visit data for 2728 Comscore panelists, showing the number of visits each person made to one website (with the disguised name of khakichinos.com). You should ignore the covariate data (we'll discuss/use it in a few weeks). Fit a plain NBD and an NBD with spike at zero to this dataset. Briefly discuss your results.

Before discussing results, we'd like to note a few design decisions for modeling for this dataset:

- For goodness-of-fit tests and model fit, we choose to censor the original data at 10 visits. This reasoning is two fold: the counts of the visits past 10 visits are low and often missing. Furthermore, the list of visits loses continuity after 22 visits. Therefore, we censor at 10 visits.
- The procedure for fitting the NBD on the full raw dataset is similar, but slightly tweaked as compared to a histogram of frequencies. For each observation, we use the gamma formula:  $P(X = x) = \frac{\Gamma(r+x)}{\Gamma(r)x!} \cdot (\frac{\alpha}{\alpha+1})^r \cdot (\frac{1}{\alpha+1})^x$  instead of forward recursion. And for the log likelihood, we simply take the logarithm of  $P(x = x)$  and sum up over all observations.
- For the goodness of fit tests, we compute probabilities according to forward recursion formulas.

## 1.1 Plain NBD

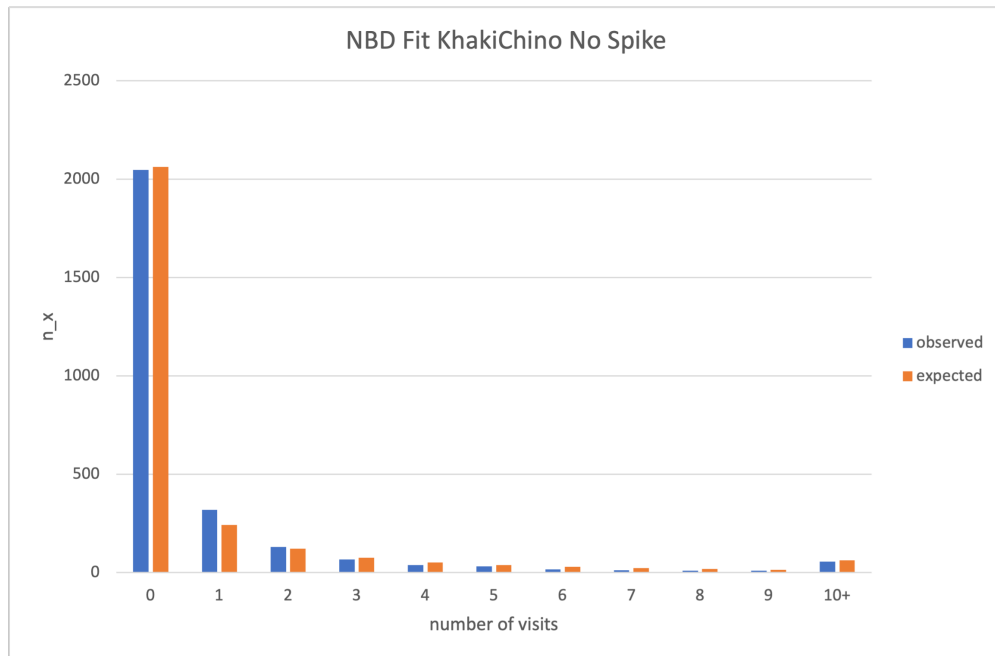


Figure 1: NBD fit of comscore data with no spike

The results of the plain nbd result in  $(r, \alpha) = (0.134, 0.141)$ . This indicates that the model believes the population of visitors to the website are heterogeneous. When fitting the chi-sq goodness-of-fit test, we use 8 degrees of freedom (11 categories, 2 parameters). The p-value for this test is  $p < 0.001$ . This indicates that the plain NBD does not explain the data very well, as the observed the data to be generated from this model would be unlikely.

Visually, the model seems to overestimate the 0 category or perhaps overestimates the number of never buyers, and tends to not fit very well across the rest of the categories.

## 1.2 NBD With Spike at 0

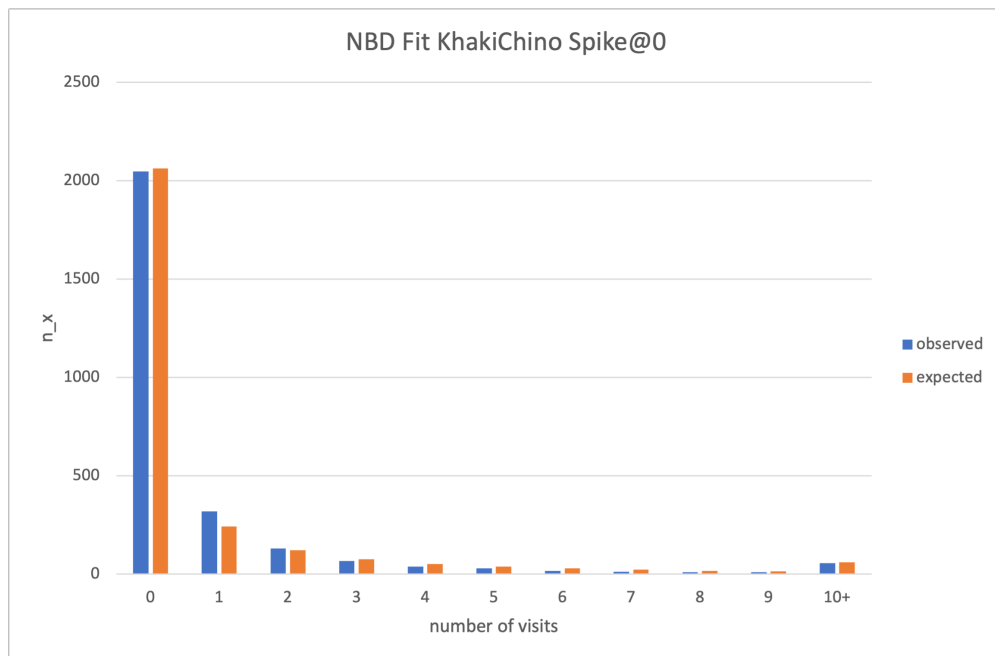


Figure 2: NBD fit of comscore data with spike@0

The results of the nbd with spike at zero result in  $(r, \alpha) = (0.134, 0.141)$  and also the spike parameter of 0.0001. This indicates that the model believes the population of visitors to the website are heterogeneous, and also that the spike seems insignificant, as a low spike parameter value indicates a low percentage of hard core never buyers. When fitting the chi-sq goodness-of-fit test, we use 7 degrees of freedom (11 categories, 3 parameters). The p-value for this test is  $p < 0.001$ . This indicates that the NBD with spike at zero does not explain the data very well, as the observed the data to be generated from this model would be unlikely.

When we compare the two models, we notice that the log likelihood that the spike model converged to is the same if not slightly higher (-2905.6266 vs -2905.6245) compared to the plain nbd model. When we compute the chi-sq LRT test, there yields no number since the difference will be negative. However, this is indicative that adding the spike parameter did not increase the performance of the model.

## 2 Question 2

The HW2A dataset contains a second sheet showing the distribution of the number of albums purchased across a set of 1574 shopping trips for an online music seller.

## 2.1 Part A: Estimate a shifted NBD and a truncated NBD. Describe the relevant features and differences between these two models.

Before we discuss the results of the modeling, we can discuss some of the implications of both of the shifted and truncated NBD models.

The shifted model assumes that potential missing categories (in this case those who never purchase albums) do not exist and we simply shift the model to starting at 1. Mechanically, this means we begin the computations of  $P(X = x)$  however many count periods later than the original model. For this model, this means that the formula for likelihood at 1 is  $(\frac{\alpha}{\alpha+1})^r$ , whereas for unshifted model this would be the probability at  $x = 0$ . We compute all the remaining probabilities as according to forward recursion formula and can perform goodness of fit in the same way. Specifically, let us define a  $\delta$  as to how many periods shifted from 0 the model becomes. Therefore, the forward recursion formula is the following:

$$P(x = x) = \begin{cases} (\frac{\alpha}{\alpha+t})^r & \text{for } x = \delta \\ \frac{t(r+x-1-\delta)}{(x-\delta)(\alpha+t)} \cdot P(x = x-1) & \text{for } x > \delta \end{cases}$$

In our case, since the shift is only one count, we have the following:

$$P(x = x) = \begin{cases} (\frac{\alpha}{\alpha+t})^r & \text{for } x = 1 \\ \frac{t(r+x-2)}{(x-1)(\alpha+t)} \cdot P(x = x-1) & \text{for } x > 1 \end{cases}$$

since the  $\delta = 1$ .

The truncated model, on the other hand, assumes the existence of missing categories, but does not fit on the data since it does not exist in the dataset. Instead, we compute probabilities for all missing categories, then normalize the final probabilities by the sum of the probabilities of missing categories. For instance, in this dataset, we normalize the rest of the probabilities by the  $1 - P(x = 0)$ , since that is missing. This also means once we fit the model, we are able to roll back the values of the missing categories to infer the counts of the missing categories. We cannot do this with the shifted model.

In the parameters  $(r, \alpha)$  of the models, we notice that there can be vast differences due to the assumption of missingness.

### 2.1.1 Shifted NBD Model

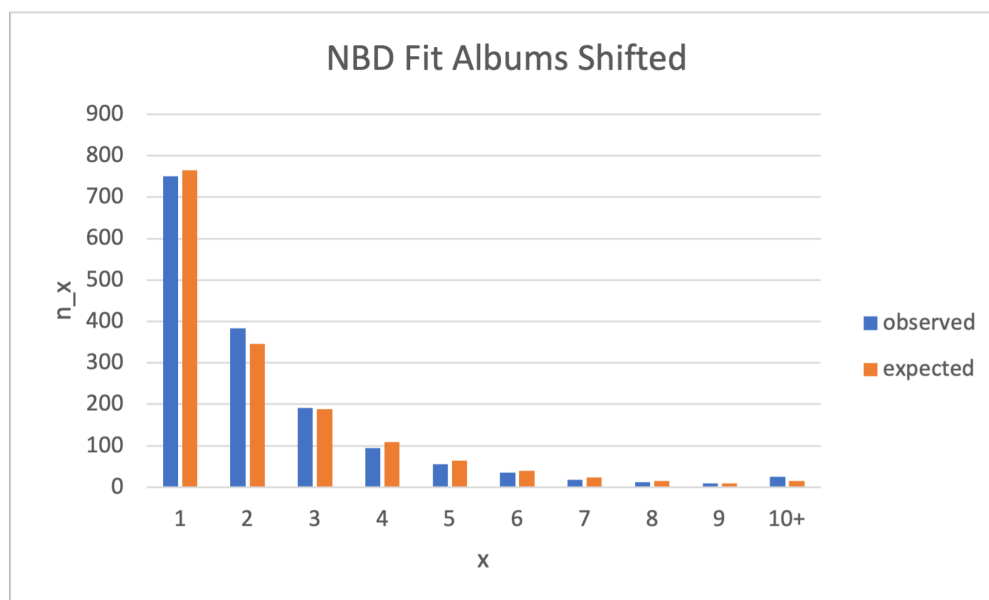


Figure 3: NBD fit of album data shifted model

The results of the shifted result in  $(r, \alpha) = (0.705, 0.560)$ . This indicates that the model believes the population of album purchasers to observe some heterogeneity. When fitting the chi-sq goodness-of-fit test, we use 7 degrees of freedom (10 categories, 2 parameters). The p-value for this test is  $p = 0.020$ . This suggests that the model, although capturing some of the shape of the data, still fails to maintain a good fit. This can be visually confirmed as we can see there are overestimates for the 1 purchase category and varying over and underestimates for the rest of the categories. Notably, the model underestimates the 10+ category (25 true vs 14.9 expected).

### 2.1.2 Truncated NBD Model

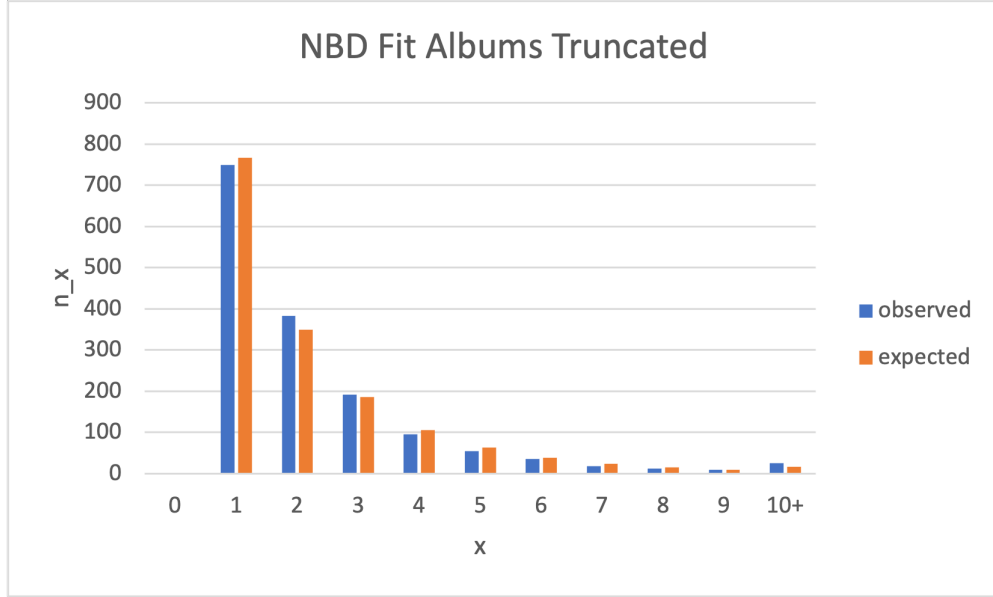


Figure 4: NBD fit of album data truncated model

The results of the truncated result in  $(r, \alpha) = (0.330, 0.461)$ . This indicates that the model believes the population of album purchasers to be fairly heterogeneous. When fitting the chi-sq goodness-of-fit test, we use 7 degrees of freedom (10 categories, 2 parameters). The p-value for this test is  $p = 0.115$ . With a p-value of 0.115, the model passes the goodness-of-fit test. This indicates that the NBD model provides a consistent fit for the data, successfully accounting for the population's heterogeneity without significant deviation from the observed purchase frequencies. Visually, the truncated model seems to fit better in most categories but seems to be missing at 1 and 2 purchase categories.

### 2.1.3 Results Comparison

Notice that the  $r$  and  $\alpha$  parameter values for the shifted and truncated models differ. The truncated model detects greater heterogeneity because it isolates the active consumer base by excluding non-purchasers. By removing these zeros, the model focuses strictly on the behavioral spread among individuals who have actually made a purchase, which typically reveals a much wider dispersion than the population-wide average suggests. Mathematically, this is reflected in the lower shape parameter  $r = 0.330$  compared to the shifted model's  $r = 0.705$ ; in the nbd, a smaller  $r$  indicates a more highly skewed Gamma distribution of purchase rates. While the shifted model effectively "smooths" the variance by forcing a start at one, the truncated model's parameters are determined solely by the differences between light and heavy buyers, thereby capturing the high degree of heterogeneity inherent in the active purchasing population.

## 2.2 Part B: Also fit and describe a “zero-truncated, one-inflated” model.

The procedure for fitting a zero-truncated, one-inflated model is simply an expansion on the truncated model. We first compute the probabilities for each category based on the truncated model and then add a spike parameter at  $x = 1$  to reflect a potential narrative for one-time buyers of albums.

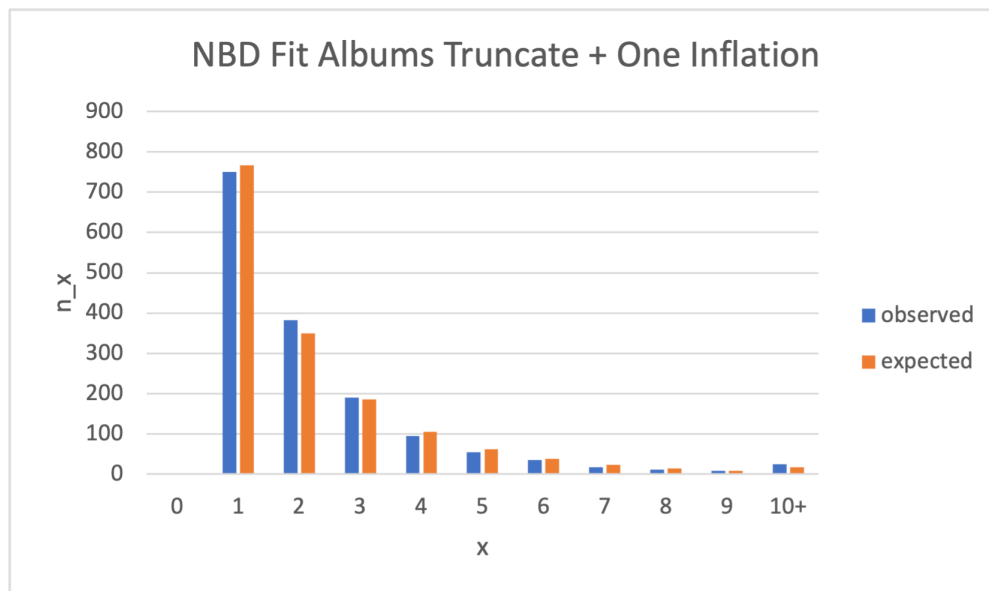


Figure 5: NBD fit of album data zero-truncated, one-inflated model

The results of the zero-truncated, one-inflated result in  $(r, \alpha) = (0.330, 0.461)$  and spike parameter of 0.00001. This indicates that the model believes the population of album purchasers to be fairly heterogeneous. When fitting the chi-sq goodness-of-fit test, we use 6 degrees of freedom (10 categories, 3 parameters). The p-value for this test is  $p = 0.070$ . The performance is similar to the original truncated model, meaning there is an absence of a one-time buyer only segment and that they do not behave differently than normal buyers of other counts.

## 2.3 Part C: Finally, fit an sBG model to the data. What would its “story” be in this context (compared to the way we used it in Session 1)?

The sBG model suggests that each shopping trip is governed by a constant probability  $p$  that the customer will stop purchasing albums in a given sitting. A shopper who purchases 9 albums is someone who ‘survived’ 8 opportunities to stop and finally decided to quit on the 9th. Crucially, the model assumes that this stopping probability varies across the 1,574 trips according to a Beta distribution, accounting for the fact that some shoppers are naturally more inclined to buy in bulk than others.

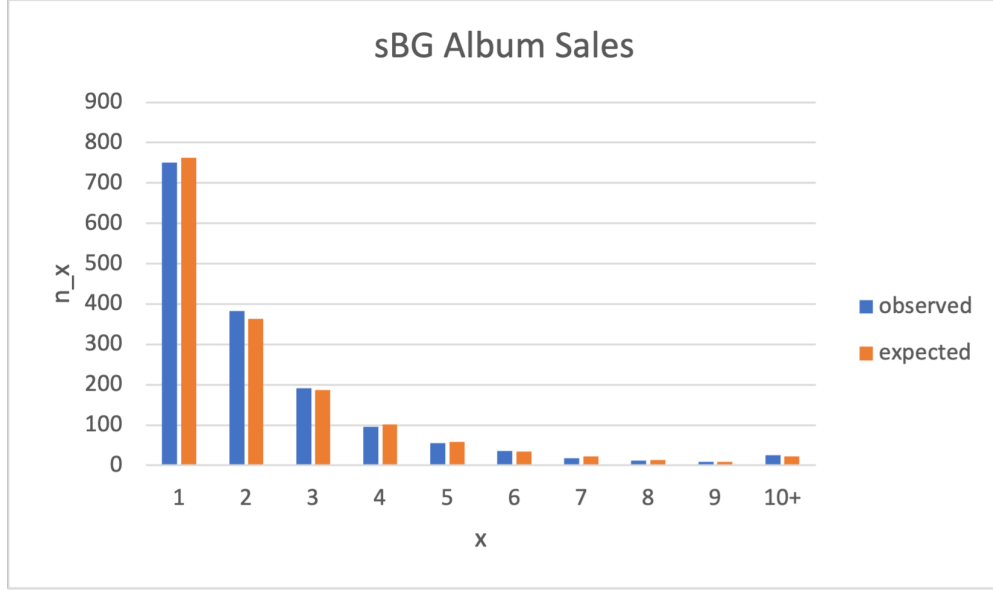


Figure 6: sBG fit of album data

The sBG model yielded parameter estimates of  $(\alpha, \beta) = (5.911, 6.282)$ . In contrast to models with parameters below 1, these relatively high and balanced values indicate a homogeneous population. Because  $\alpha$  and  $\beta$  are similar and significantly greater than 1, the underlying Beta distribution of  $p$  (the stopping probability) is unimodal and roughly symmetric, centered near 0.5. This "bell-shaped" distribution suggests that most shoppers share a similar, moderate propensity to stop after each purchase, rather than being split into extreme "power buyers" and "one-time" shoppers. For the Chi-squared goodness-of-fit test, we utilized 7 degrees of freedom (10 categories, 2 parameters). The resulting p-value of  $p = 0.855$  is exceptionally high, suggesting that the model provides an excellent fit for the observed data.

## 2.4 Part D: Which of these 4 models would you choose and why?

To make a decision for the best model, we compare each based on their statistical fit and the underlying "story" they tell about consumer behavior. While the sBG model provides a statistically superior fit in this instance with a p-value of 0.855, its narrative—that consumers engage in a sequential "coin-flip" process where they decide to stop after each individual album—can feel implausible in a real-world retail setting. In contrast, we often appreciate models that view behavior as counts over a time period, such as the NBD model. This perspective suggests that consumers possess a latent, long-term "purchase rate," which is a more intuitive reflection of a person's stable interest in music. Rather than a series of impulsive stopping decisions, the NBD narrative allows us to see a shopper who buys ten albums not as someone who "survived" nine trials, but as a "heavy buyer" with a high consumption velocity. Ultimately, while the sBG captures the mechanics of a single shopping trip with high precision, the NBD provides a more managerially useful story by linking observed counts to the enduring traits of the customer.





NBD Fit KhakiChino Spike@0

$n_x$	observed	expected
0	2050	2080

Figure 8: nbd model on comscore data spike@0

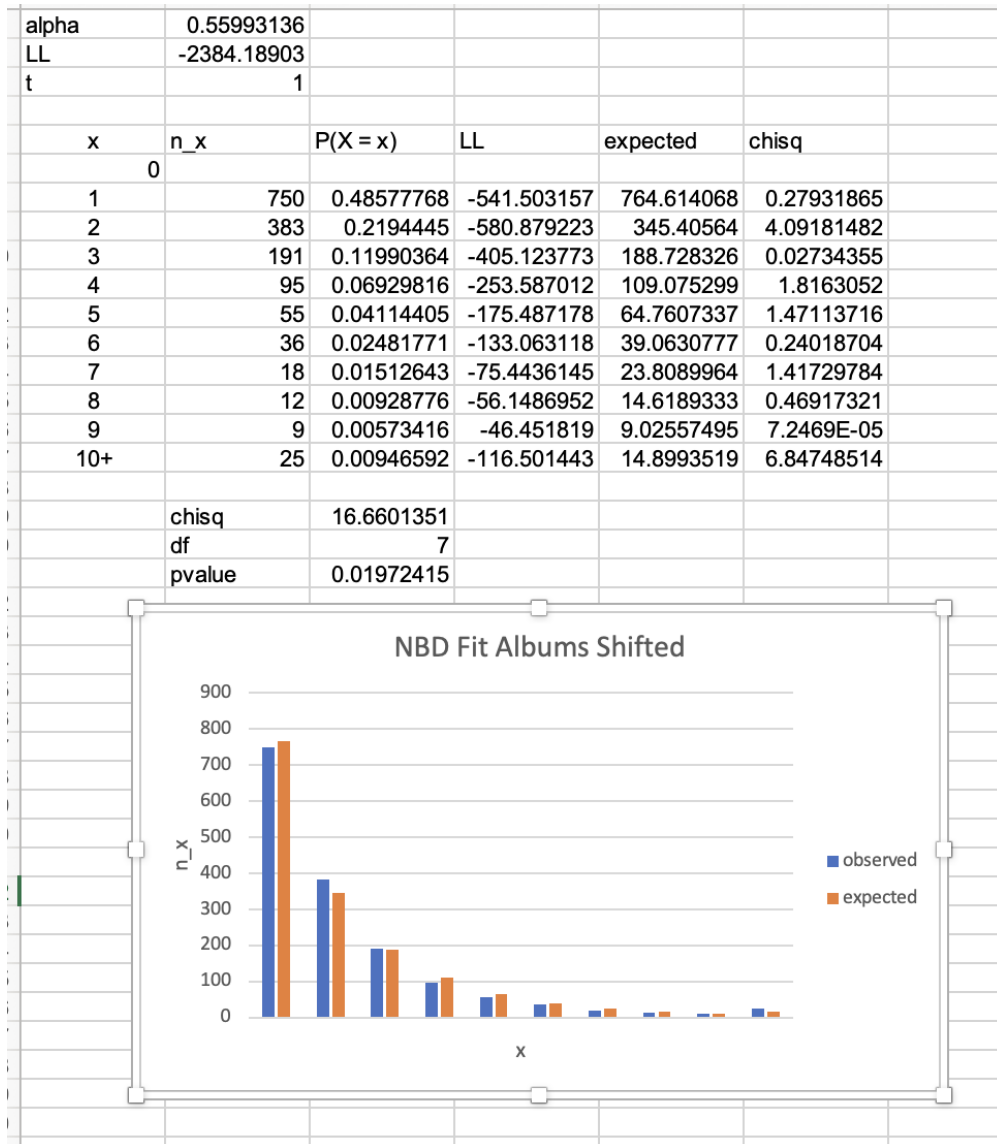


Figure 9: shifted nbd model on albums

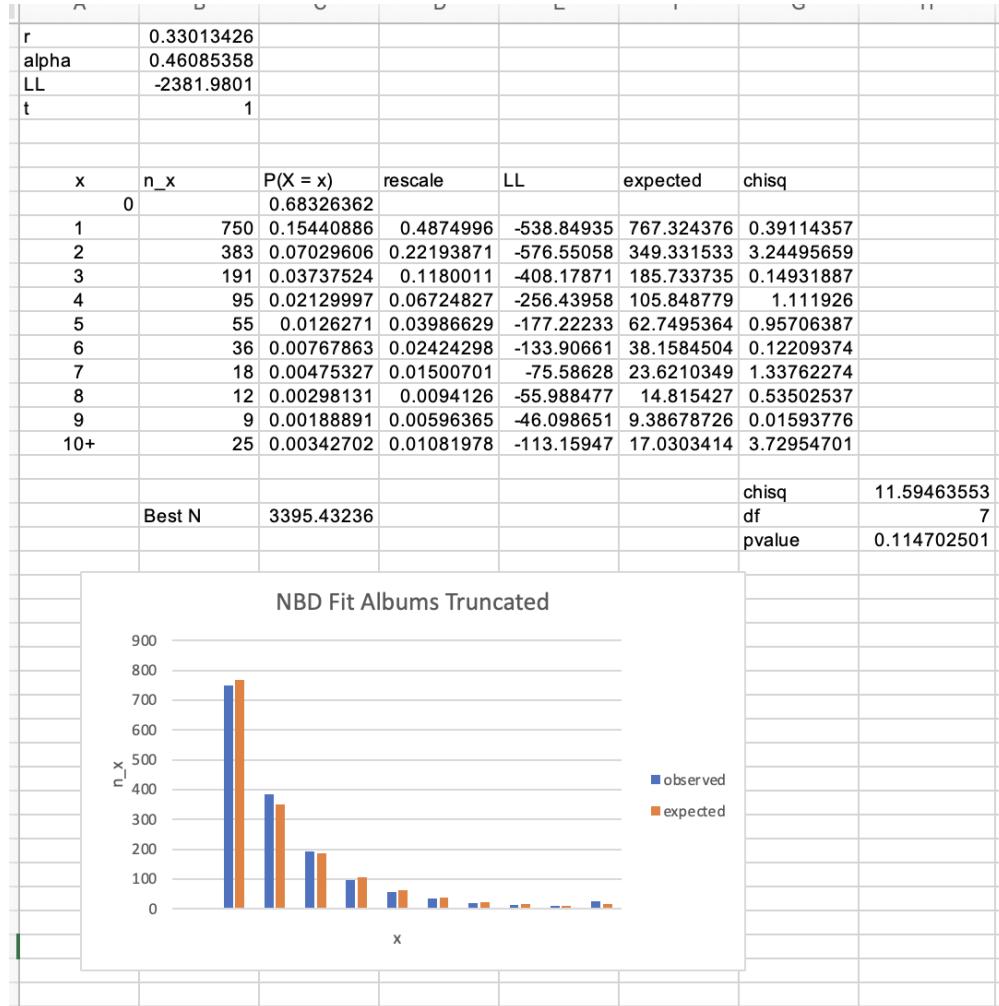


Figure 10: truncated nbd model on albums

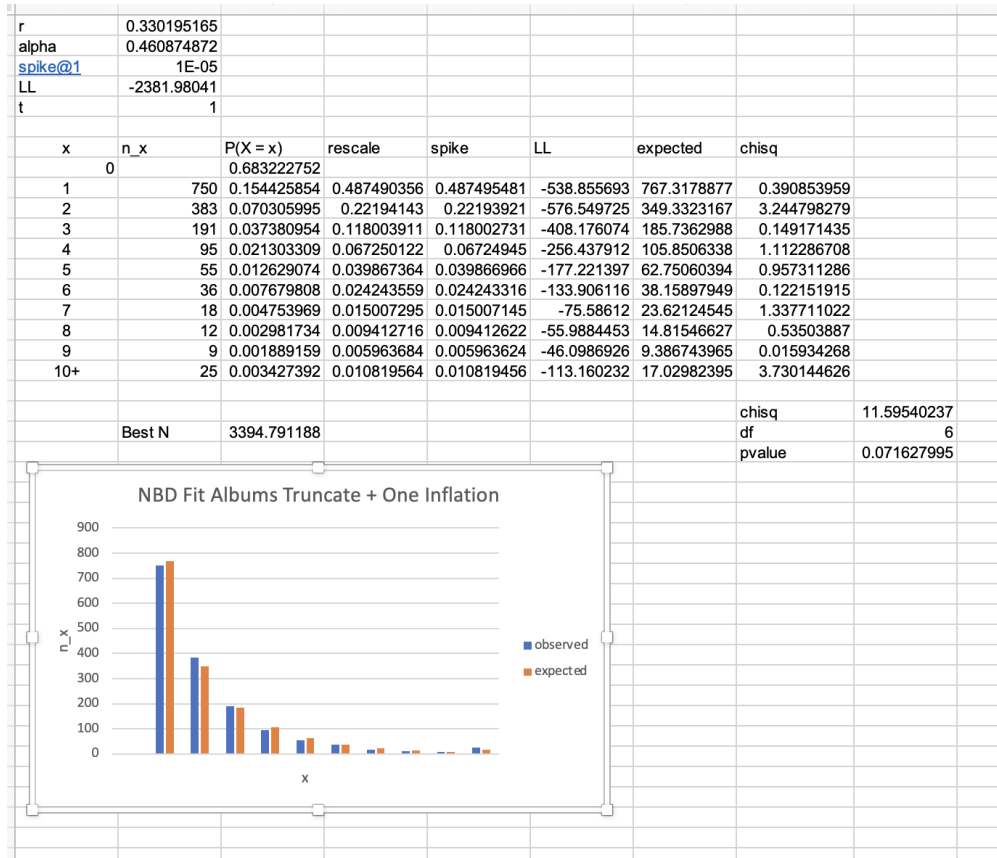


Figure 11: zero-truncated, one inflated nbd model on albums

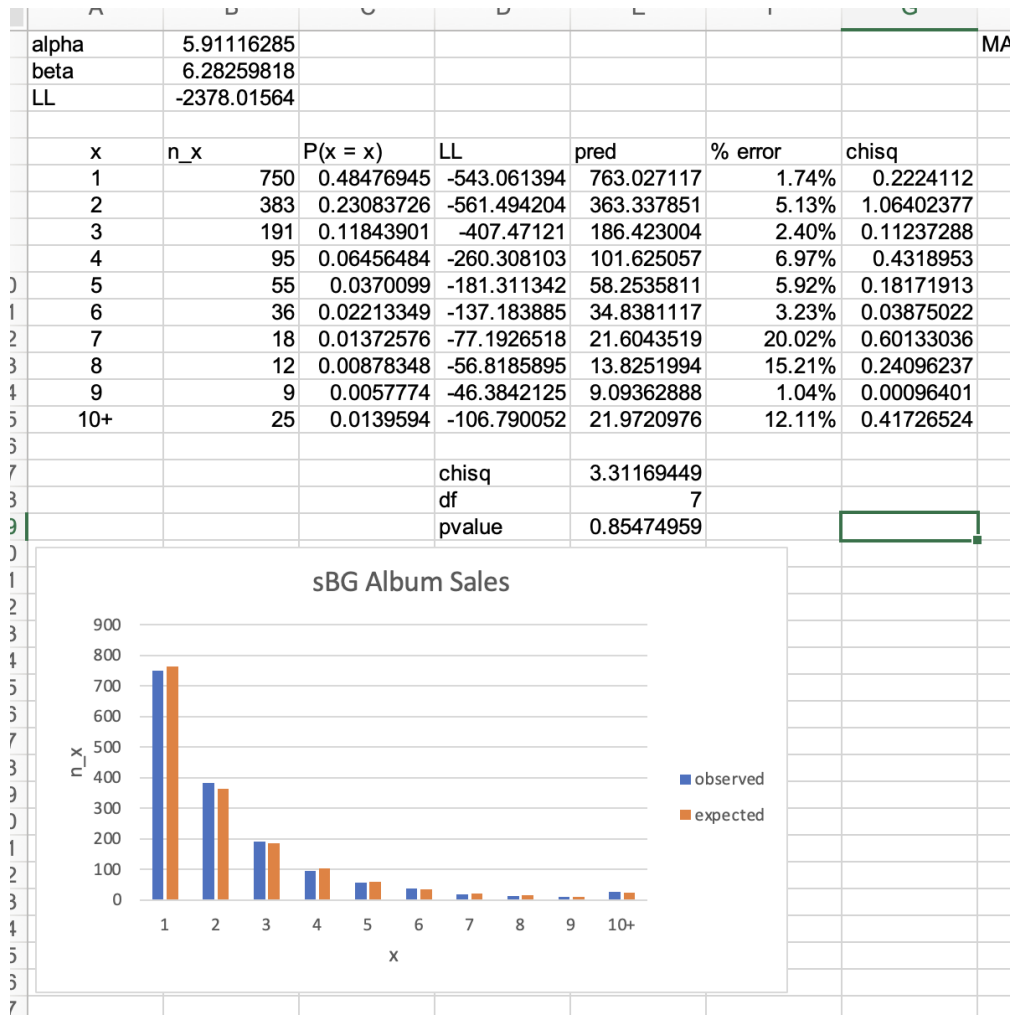


Figure 12: sbg model on albums