

Gravity Model

Twitter Data from NYC

Overview

1. Probabilistic Weight

Construct a network with the probabilistic weight of link(a,b) defined as

$$\text{link weight}(a, b) = \sum_c \frac{t(c, a) \cdot (t(c, b) - \delta(a, b))}{T \cdot (t(c) - 1)}$$

where

- $t(c, a)$ denotes the total number of tweets that user c has posted at location (in our case, zipcode) a
- $t(c) = \sum_a t(c, a)$
- $T = \sum_c t(c)$

2. Gravity Model

On the other hand, the weight of links can be modeled as

$$\text{link weight}(a, b) = k \cdot w^{\text{out}}(a) \cdot w^{\text{in}}(b) \cdot f(d(a, b)),$$

where

- $w(x)$ is the weight or customized centrality of the node x
- $f(d(a, b))$ denotes function with respect to distance between location a and b , usually it decays as distance increases.

Formulae

So after calculating all $link(a, b)$, we run linear regression by taking logarithm on both sides

A few notes:

- **Original model** uses link weight directly from the formula:

$$\log(weight_{a,b}) \sim \log(w_a^{out}) + \log(w_b^{in}) + \log(f(d(a, b)))$$

- **Adjusted model** uses (link weight) * (number of tweets from origin) instead:

$$\log(weight_{a,b}) + \log(number\ of\ tweets\ at\ a) \sim \log(w_a^{out}) + \log(w_b^{in}) + \log(f(d(a, b)))$$

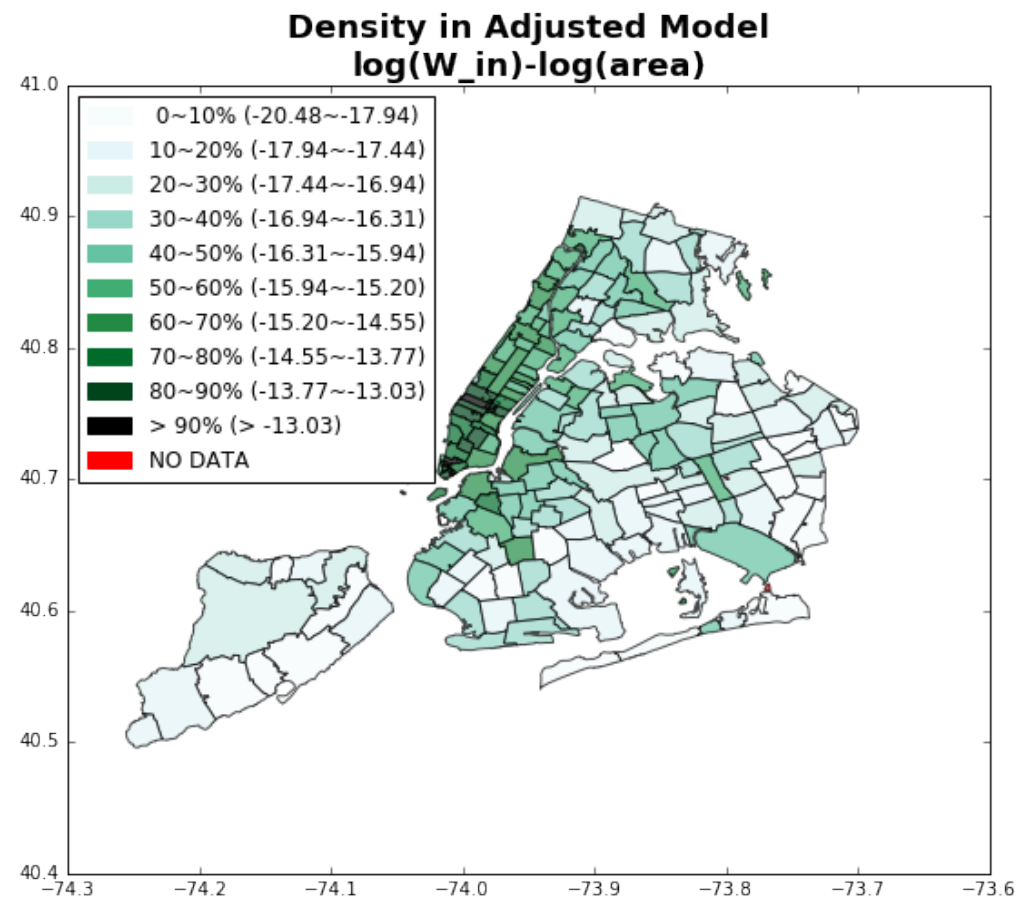
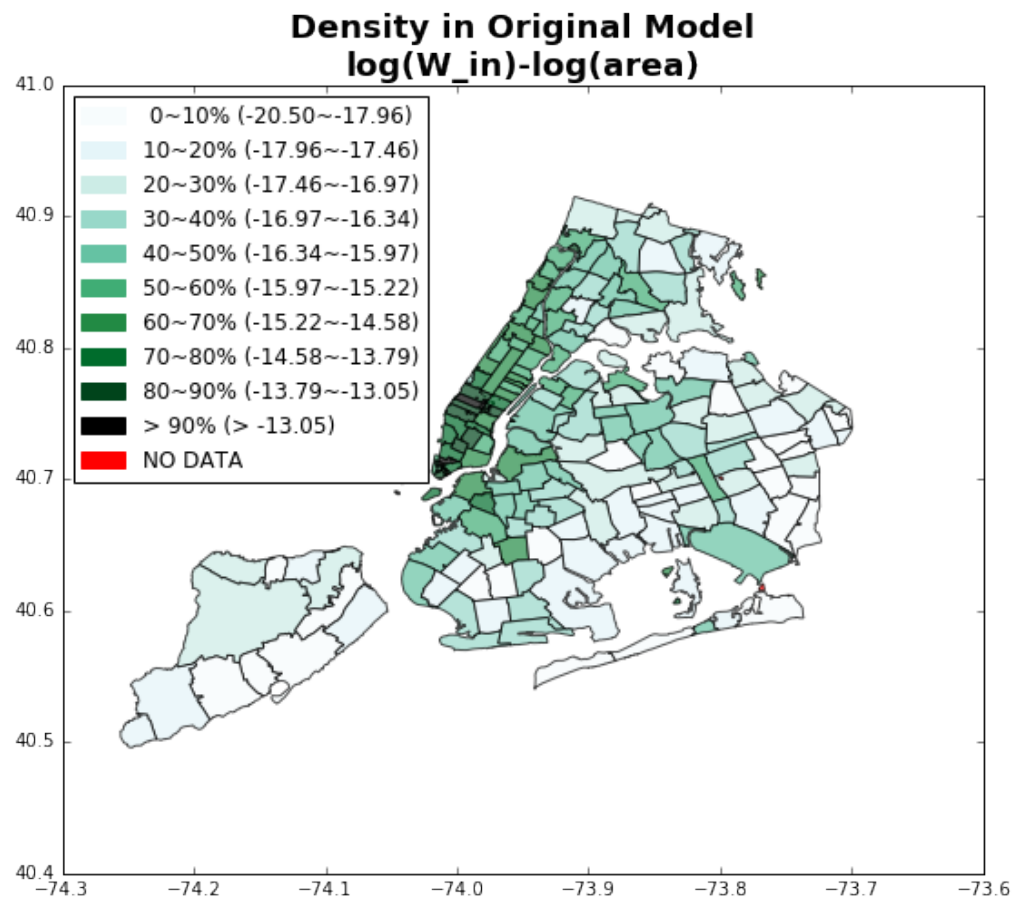
- NYC has about 240~250 zip codes, so we have 242 terms each(from twitter data) to be fit, for both W^{out} and W^{in}
- For $f(d(a, b))$, we used a custom binned function to separate all distance data into 100 subcategories based on its position between percentiles, so we have to fit an additional amount of 100 of distance bin, then observe the relationship between the coefficients and the distance within each bin

On the following pages, we use heatmap to showcase both models for comparison. As we can see from the plot, adding $\log(number\ of\ tweets\ at\ a)$ does significantly change the look of W_a^{out} , but W_b^{in} appears little difference

W_b^{in} :

Divide by the area

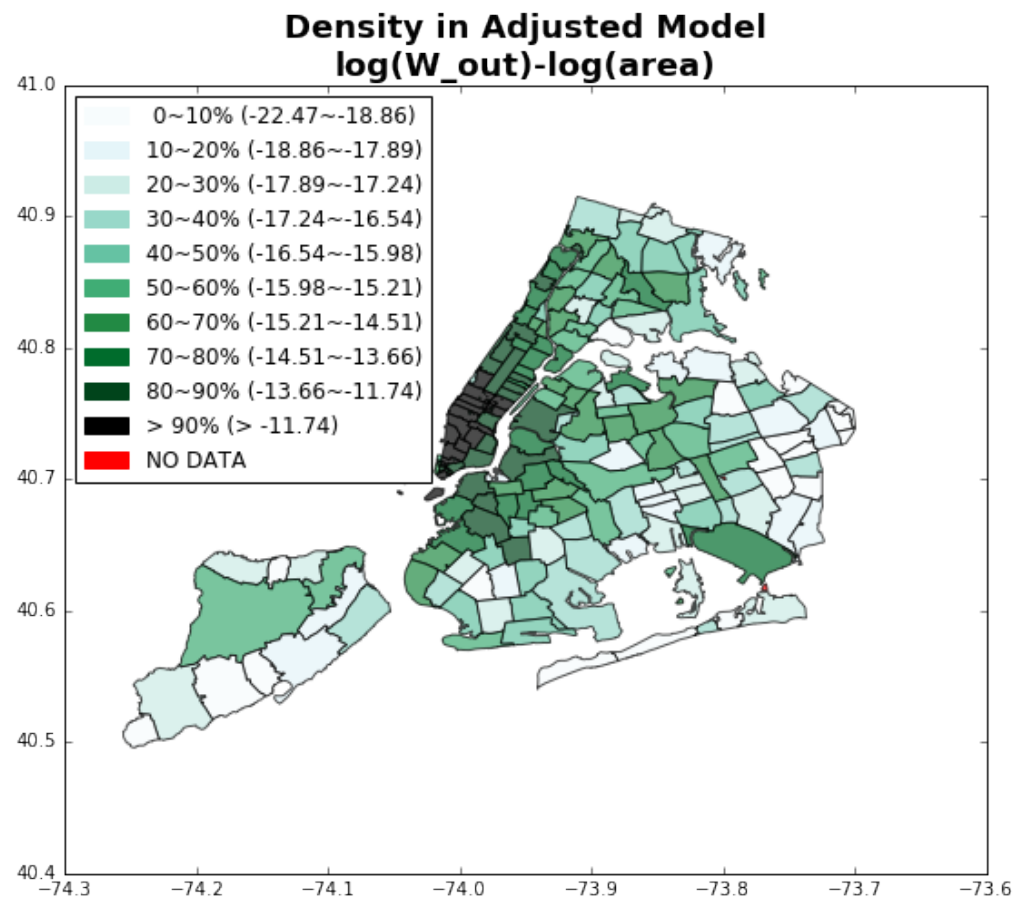
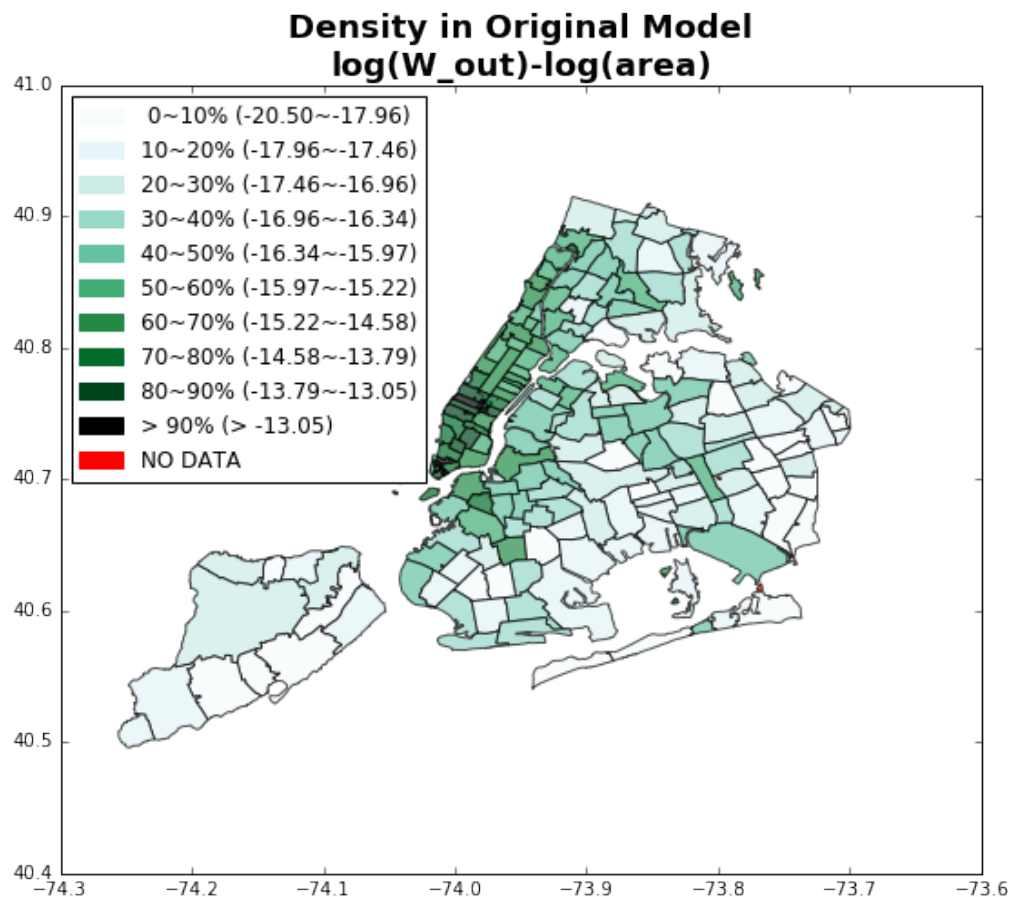
Not showing much difference if we add the log of number of tweets



W_a^{out} :

Divide by the area as well

The absolute quantities don't change significantly but the distribution does, resulting in more darker coverage all over the city on different level



Further exploration

We've focused on the dataset before Dec. 19, 2015 to avoid massive empty data frames between Dec. 19, 2015 and Feb. 1, 2016

- 1. Drop constant term in OLS, see what will come out
- 2. Examine distance function model

$$\log(f(x)) = a \cdot \log(x) + b$$

or equivalently

$$f(x) = k \cdot x^a$$

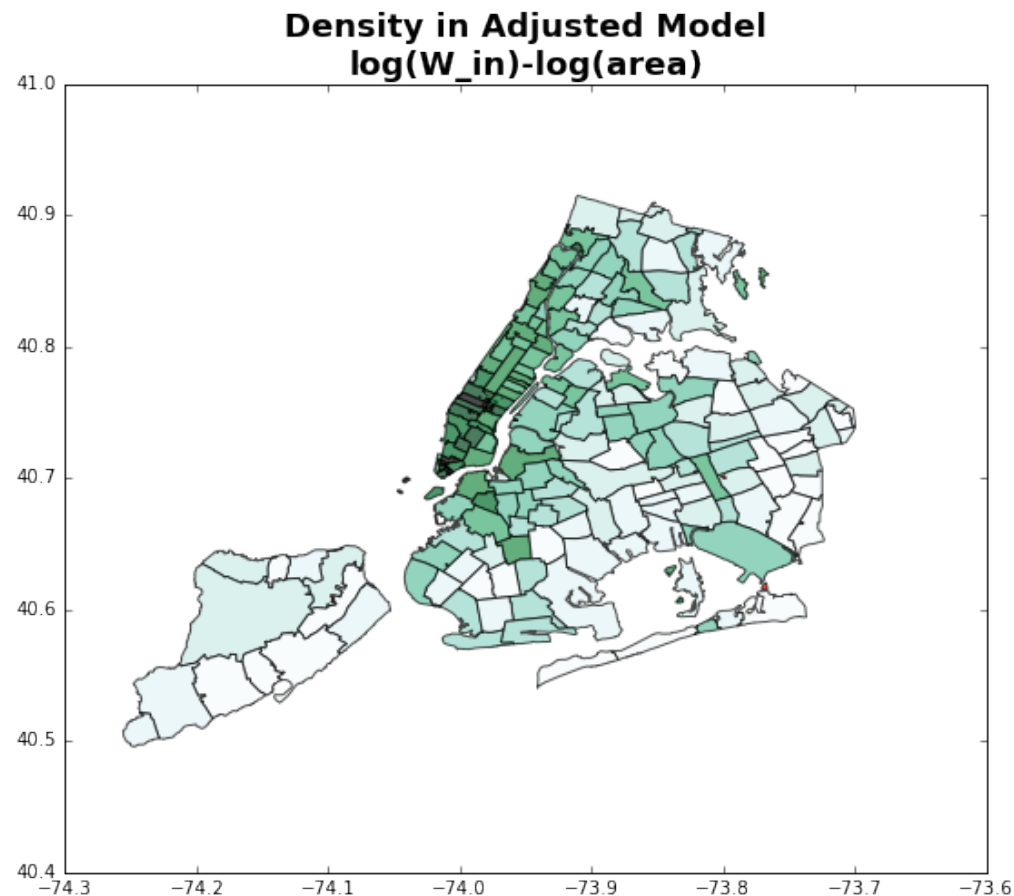
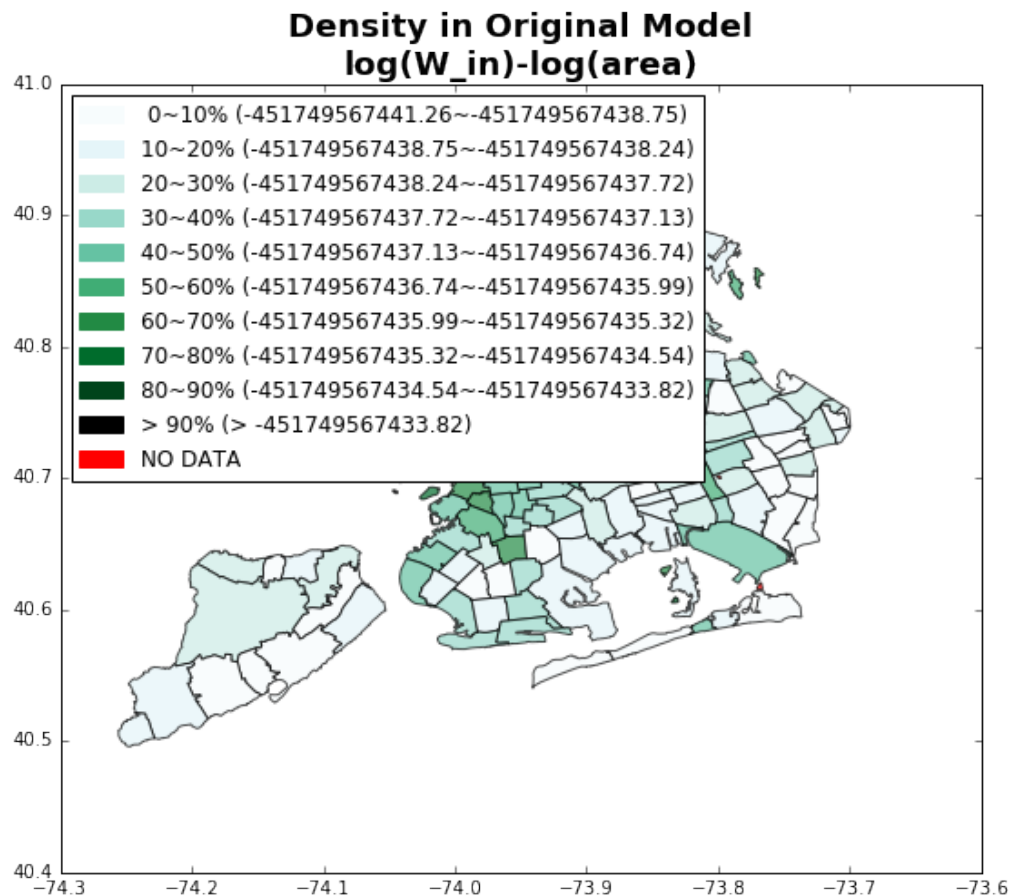
where

- $\log(f(x))$ is the fitted coefficients for each distance bin in our previous *adjusted model*
 - x denotes the distance (we've converted the average distance within each bin to this variable)
 - $b = \log(k)$ denotes the intercept in OLS
- 3. Check the coefficient a and b for each day, see if they vary drastically, and compare it with the overall fitting number.

Drop the constant

Calling `result2 = ols(y=Y2, x=X, intercept = False)`

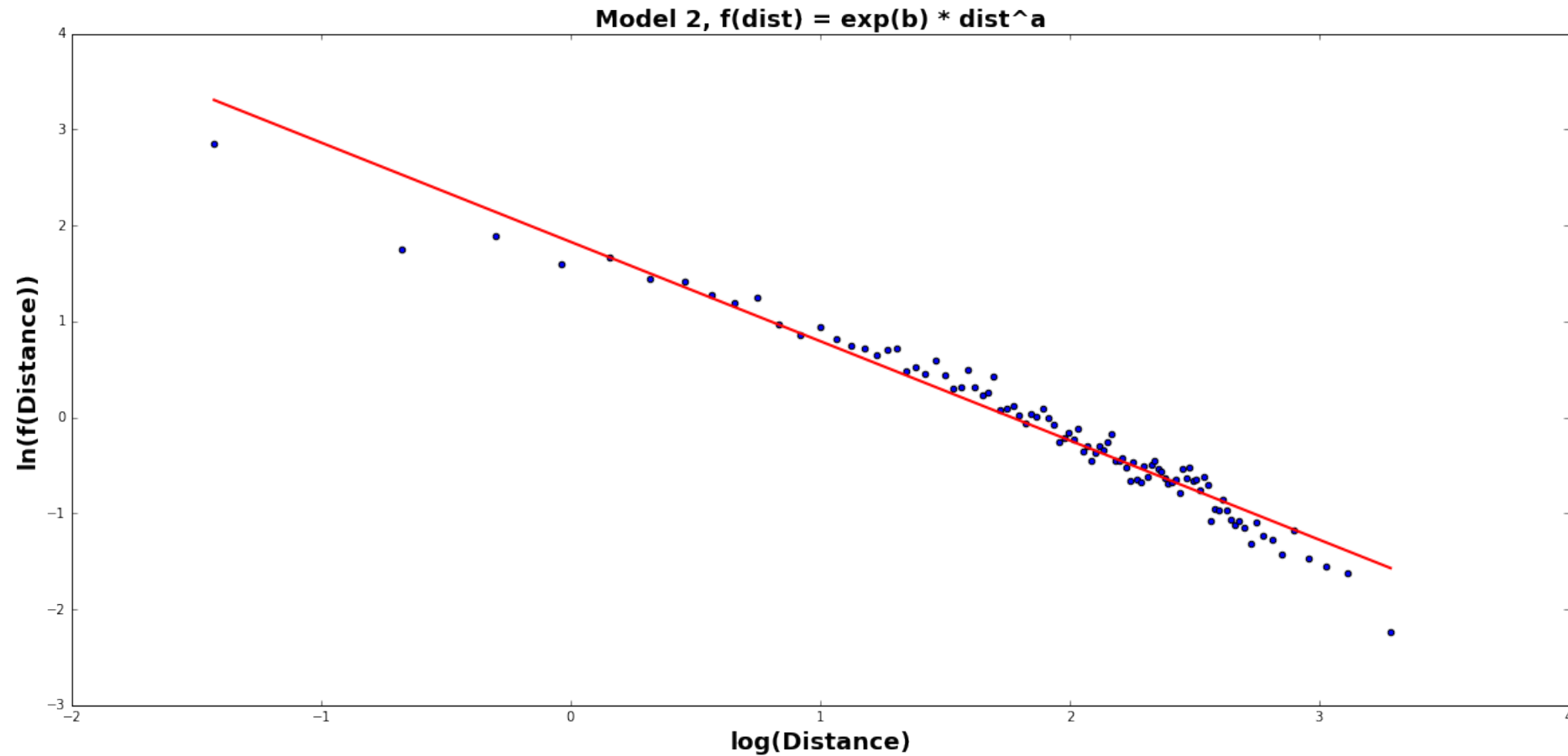
all coefficients will go crazy (left figure) but the color won't change much (right figure)



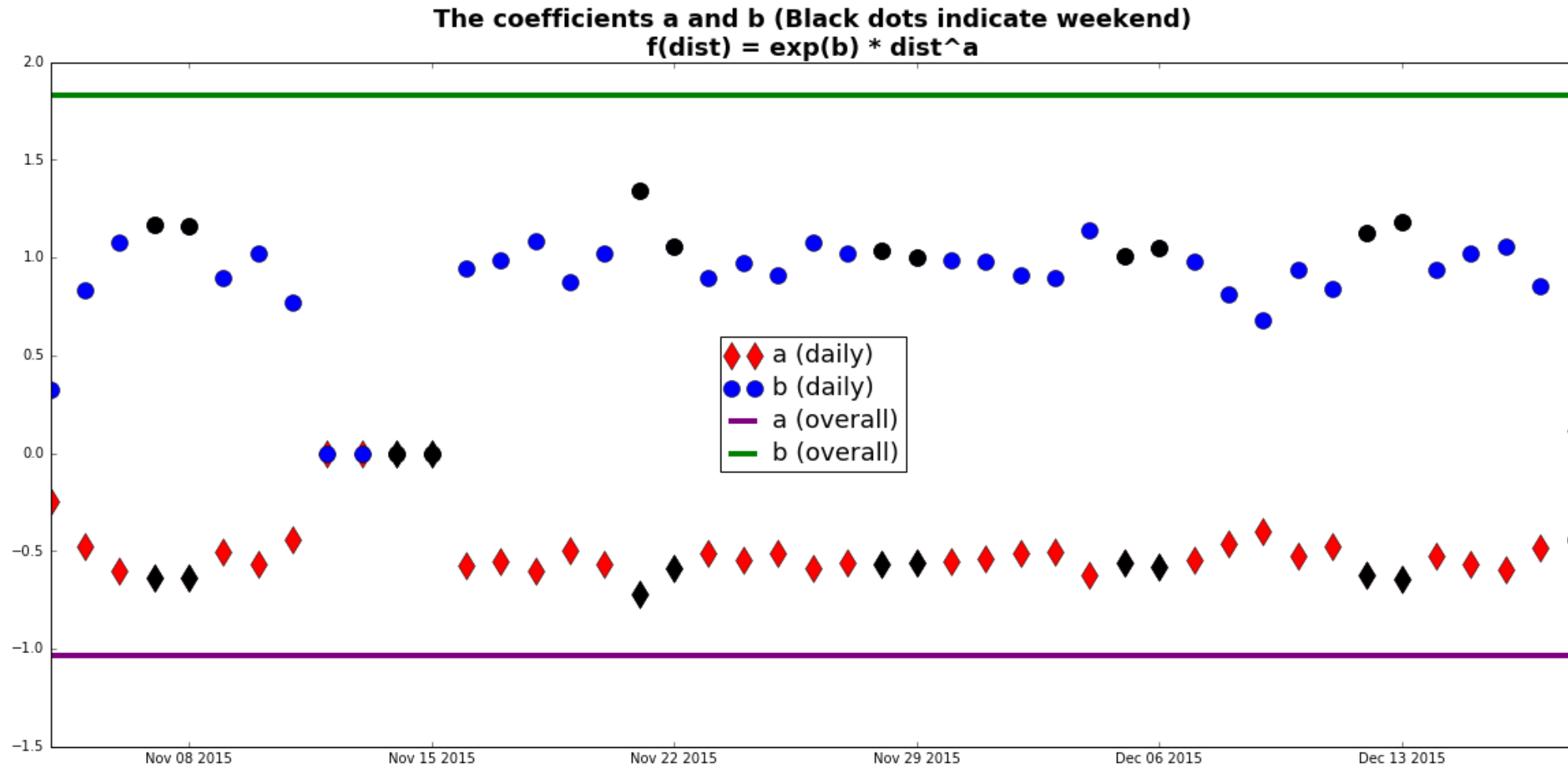
Overall estimate for a and b

The coefficients of this model are $a = -1.034442$, $b = 1.828255$

The R square of this model is 0.958063



Observe a and b on each day



...daily a and b continued

- For non-zero part, standard deviations are:
 - Coef_A 0.128473
 - Coef_B 0.273142
- For non-zero part, average values are:
 - Coef_A -0.528674
 - Coef_B 0.937144
- Compared to the overall values of a and b:
 - Coef_A -1.034442
 - Coef_B 1.828255