



What Makes for a Popular Book Review?

Michael Burnham-Fink
Metis Project Luther

Business Case

- Popular book reviews are engaging

good**reads**



METIS

Business Case

- Engagement on Goodreads.com impacts book sales on Amazon

The Amazon logo is centered within a large, light gray circle. The logo itself consists of the word "amazon" in a bold, black, sans-serif font, with a curved orange arrow underneath it pointing from the 'a' to the 'z'.

METIS

Tools

Scraping



Modelling

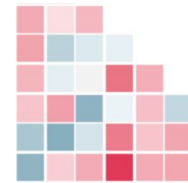


StatsModels

Visualizing



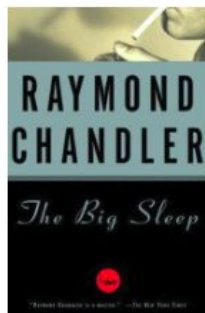
BeautifulSoup



seaborn

Data Source

Michael Burnam-Fink's Reviews > The Big Sleep



The Big Sleep (Philip Marlowe, #1)

by Raymond Chandler



Michael Burnam-Fink's review

May 05, 2012 · edit



bookshelves: mystery, fiction, 2012

7 highlights (Private)

What can you say about "The Big Sleep" that hasn't already been said? This is a classic of hard boiled noir. The language is as glamorous as the dames, as punchy as private-eye Marlowe, and as dark as the souls of the criminals, dissolute rich, and corrupt cops who inhabit the world. If you haven't read it, you're missing out.

1 like

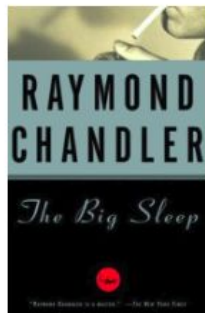
READING PROGRESS

- Started Reading [Add a date](#)
- May 5, 2012 - Shelved
- May 5, 2012 - Shelved as: [mystery](#)
- May 5, 2012 - Shelved as: [fiction](#)
- May 5, 2012 - Shelved as: [2012](#)
- May 5, 2012 - Finished Reading

Data Source

- Word Count
- Image Count
- Book info
- Rating
- User Friends
- Average Rating
- Total Times Reviewed

Michael Burnam-Fink's Reviews > The Big Sleep



The Big Sleep (Philip Marlowe, #1)

by Raymond Chandler



Michael Burnam-Fink's review

May 05, 2012 · edit



bookshelves: mystery, fiction, 2012

7 highlights (Private)

What can you say about "The Big Sleep" that hasn't already been said? This is a classic of hard boiled noir. The language is as glamorous as the dames, as punchy as private-eye Marlowe, and as dark as the souls of the criminals, dissolute rich, and corrupt cops who inhabit the world. If you haven't read it, you're missing out.

1 like

Target

READING PROGRESS

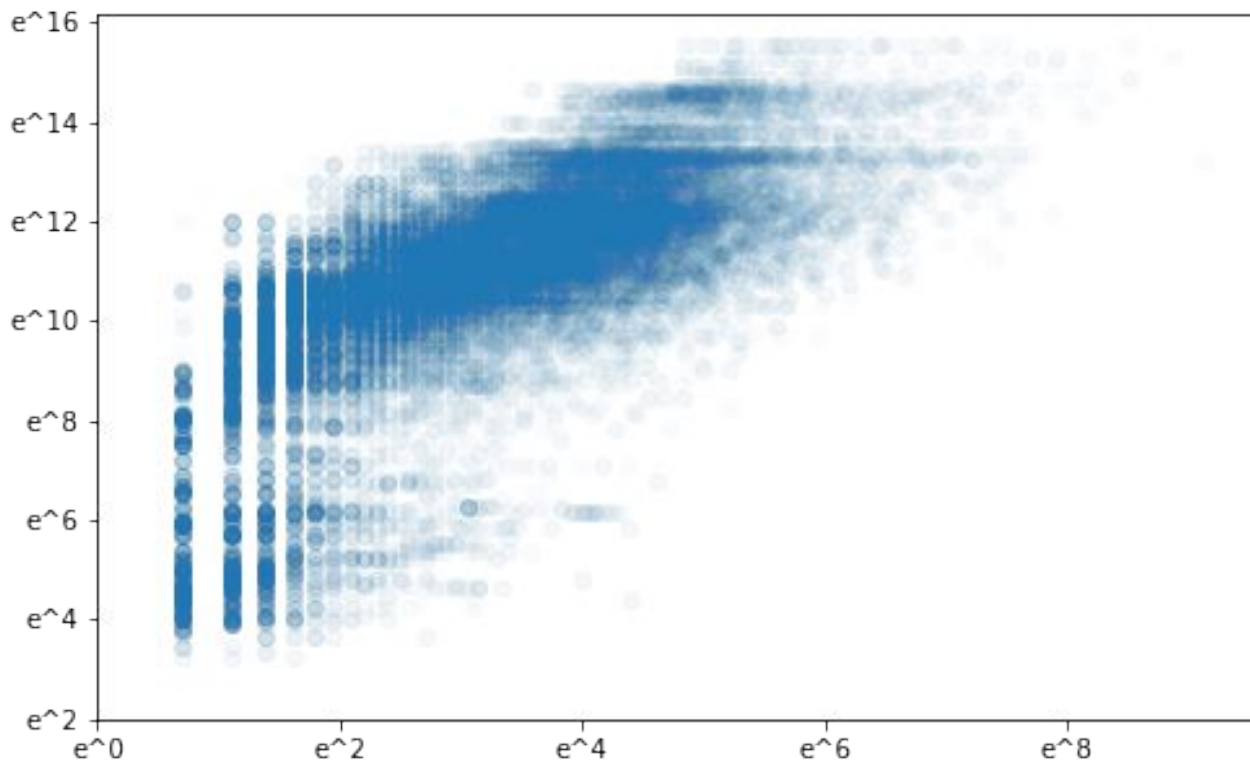
- Started Reading [Add a date](#)
- May 5, 2012 - Shelved
- May 5, 2012 - Shelved as: [mystery](#)
- May 5, 2012 - Shelved as: [fiction](#)
- May 5, 2012 - Shelved as: [2012](#)
- May 5, 2012 - Finished Reading



METIS

EDA

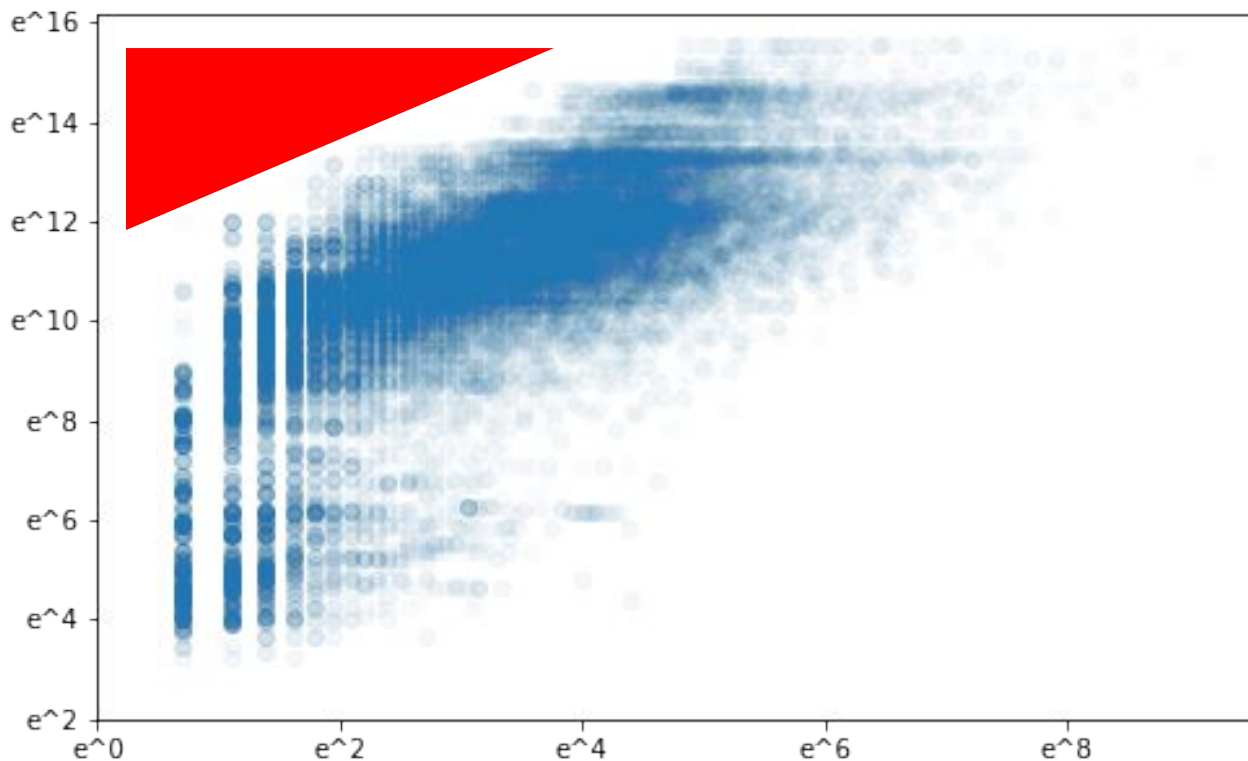
Relationship between 'Review Likes' and 'Times Book Rated'



EDA

**THERE
SHOULD
BE DATA
HERE**

Relationship between 'Review Likes'
and 'Times Book Rated'

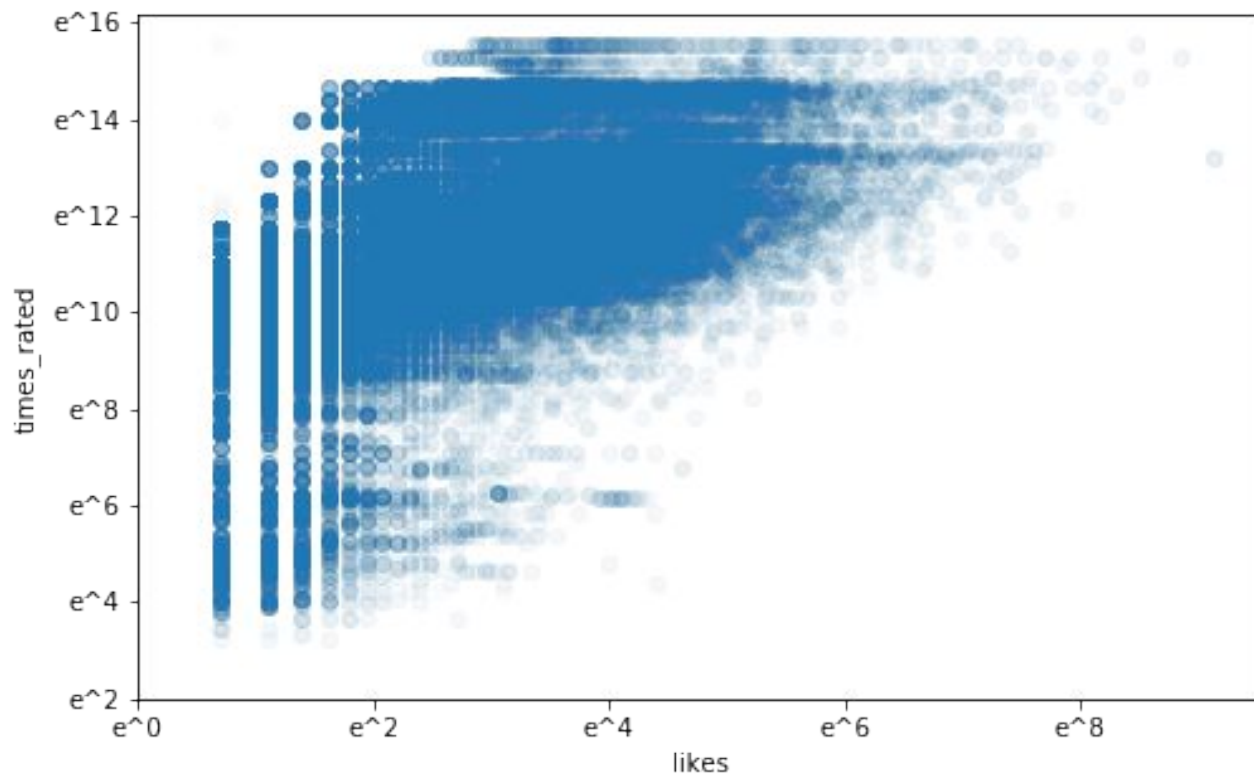




EDA

75,000 records
200 MB of data

Relationship between 'Review Likes' and 'Times Book Rated'



Code: Multiprocessing

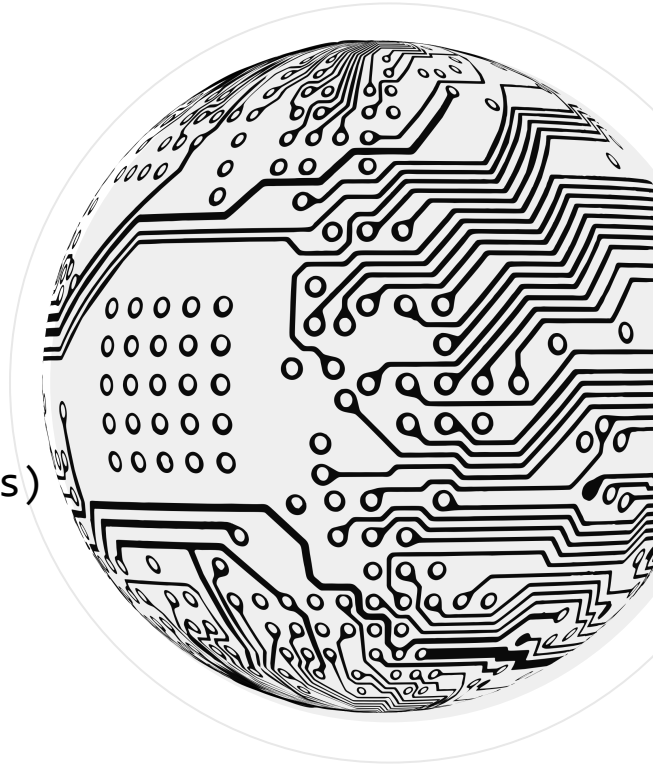
```
import multiprocessing

def reviewscraper(url)
review_urls = [url_1, url_2... url_300]

pool = multiprocessing.Pool()
reviews = pool.map(review_scraper, review_urls)

df = pd.DataFrame.from_records(reviews)

outfile = open('title book.pkl','wb')
pickle.dump(df,outfile)
```

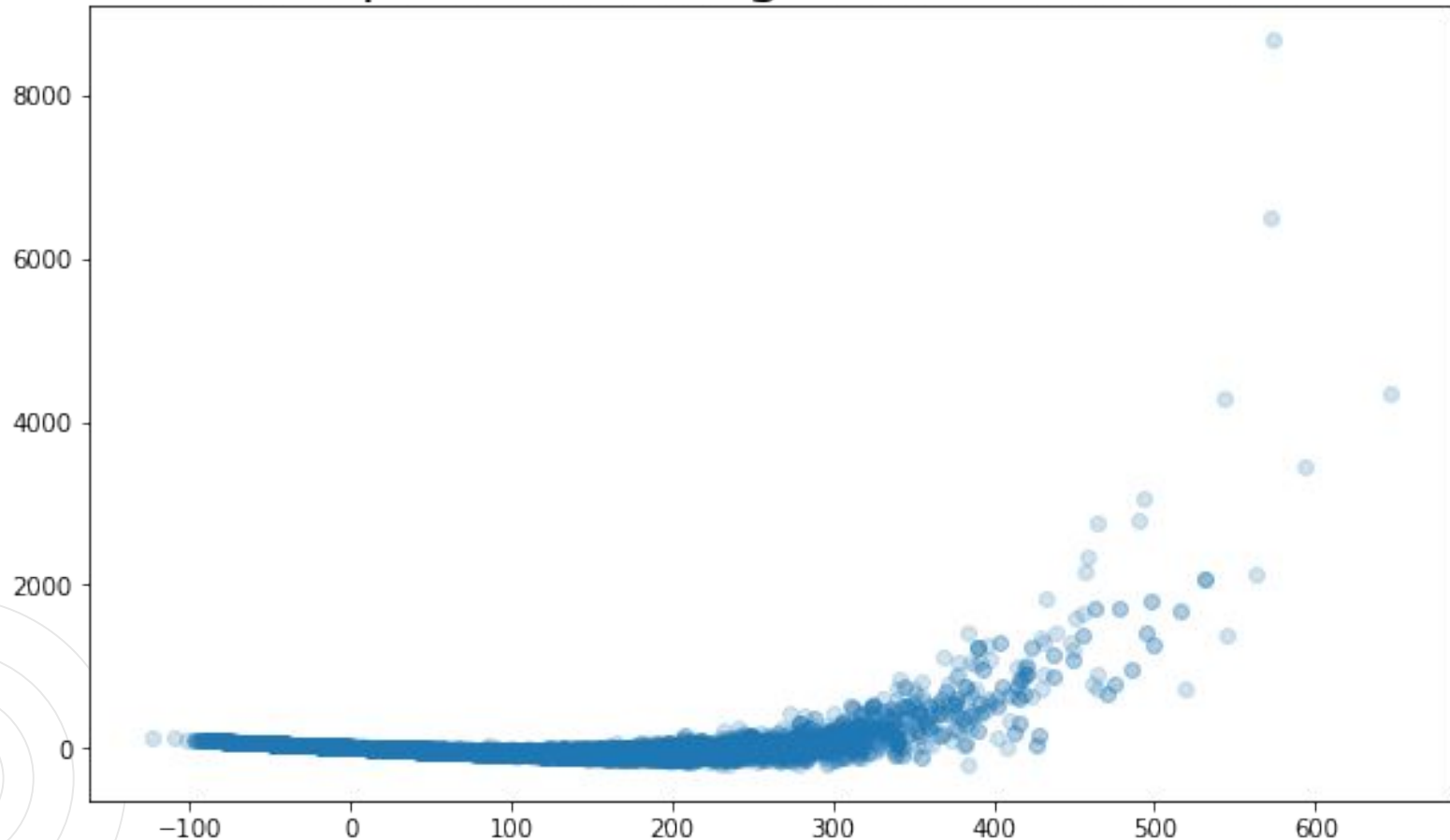


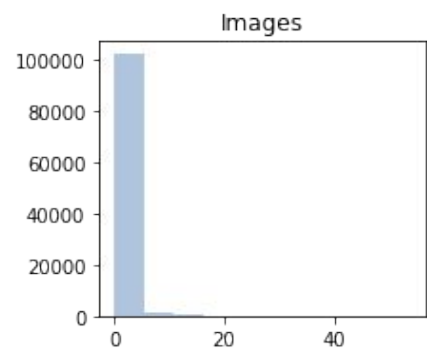
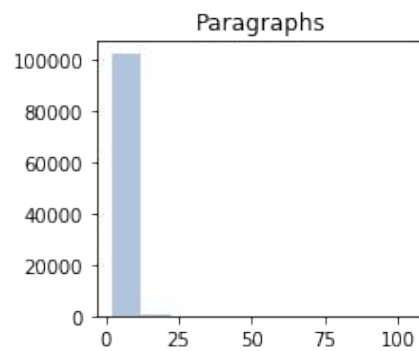
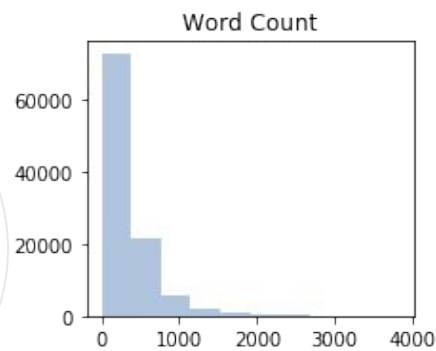
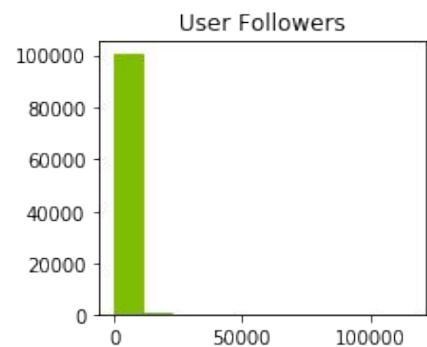
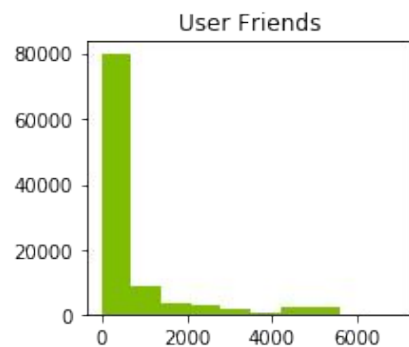
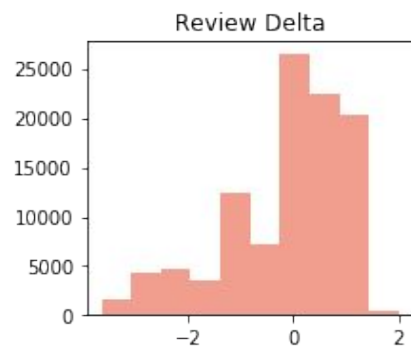
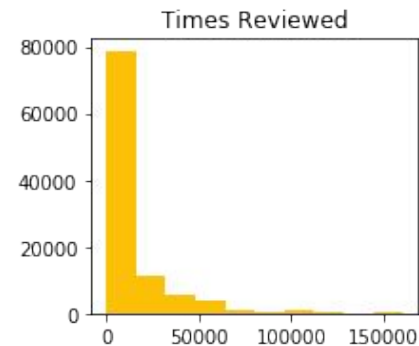
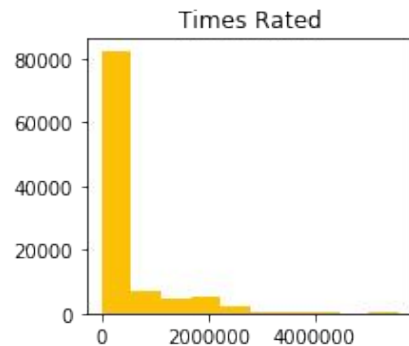
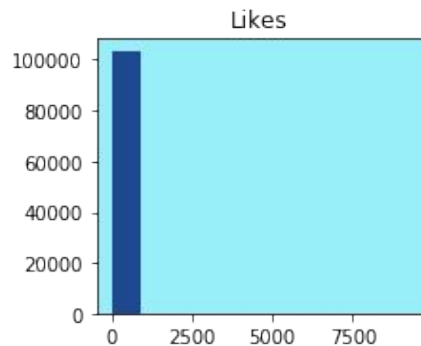
First Linear Regression


Dep. Variable:	y	R-squared:	0.256
Model:	OLS	Adj. R-squared:	0.255
Method:	Least Squares	F-statistic:	1226.
Date:	Thu, 11 Oct 2018	Prob (F-statistic):	0.00
Time:	14:28:31	Log-Likelihood:	-3.1845e+05
No. Observations:	53576	AIC:	6.369e+05
Df Residuals:	53560	BIC:	6.371e+05
Df Model:	15		
Covariance Type:	nonrobust		

R^2 : 0.256
MSE: 5,883

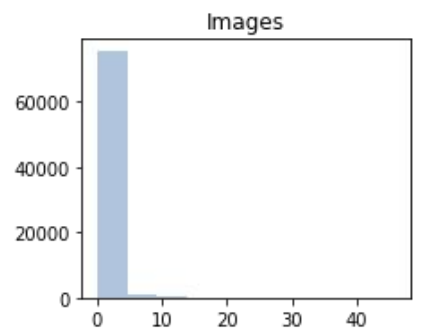
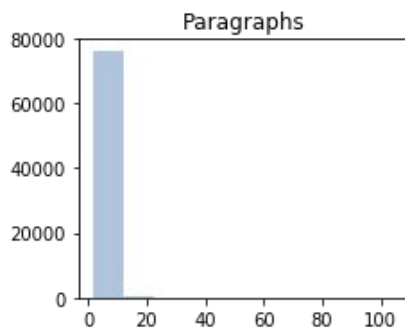
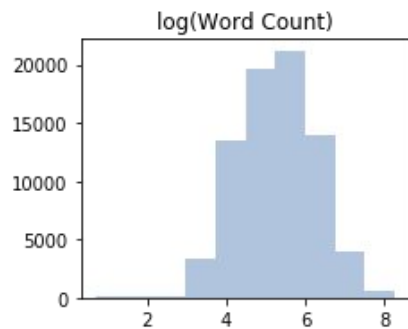
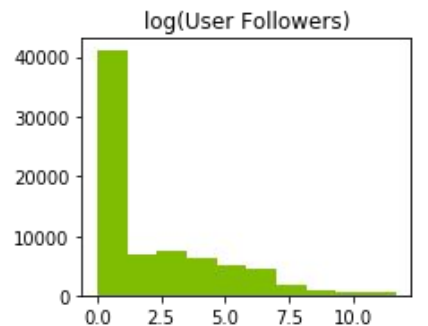
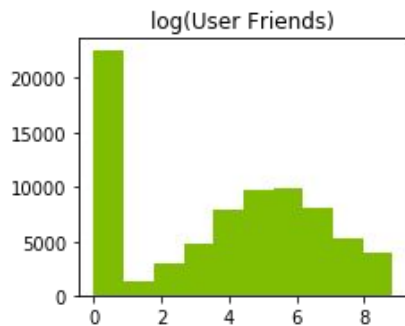
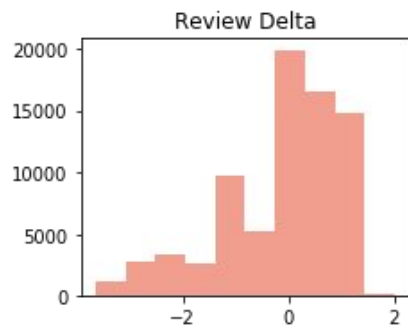
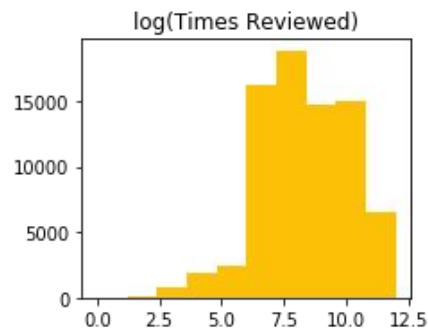
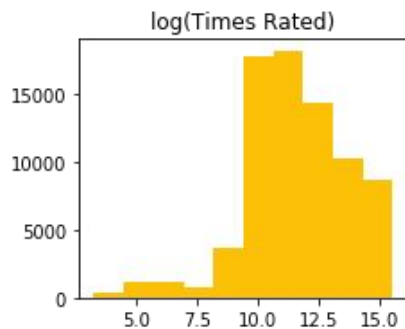
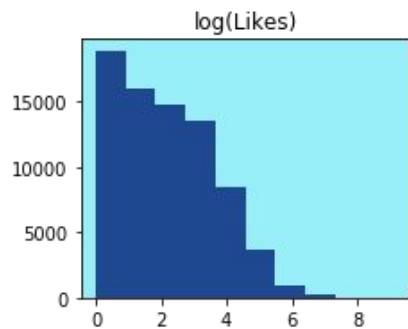
Simple Linear Regression Residuals



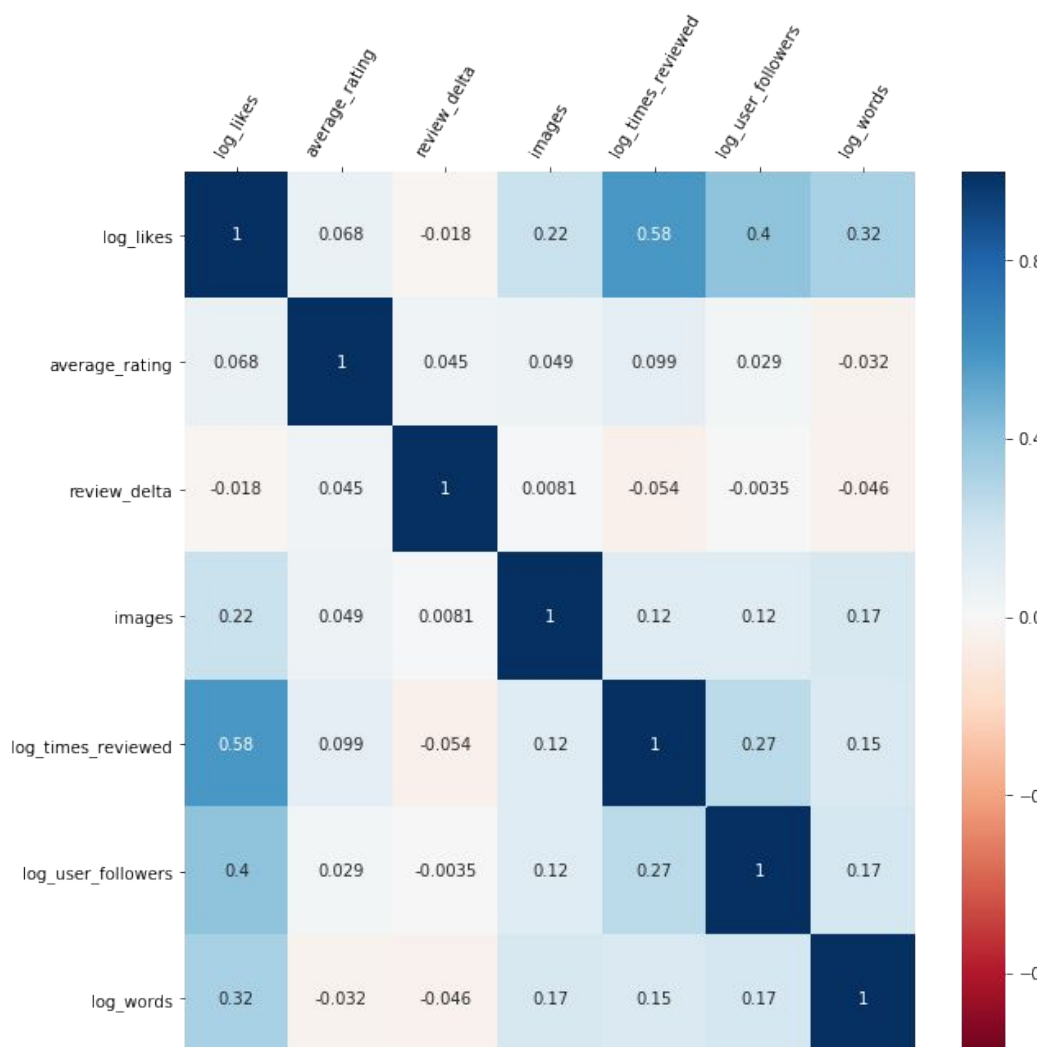


The background of the slide features a series of concentric circles in a light gray color, centered on a white background. The circles vary in size, creating a subtle, modern design.

Social networks
grow by
preferential attachment
and demonstrate
exponential distributions



METIS



Final LR

Dep. Variable:	log_likes	R-squared:	0.461
Model:	OLS	Adj. R-squared:	0.461
Method:	Least Squares	F-statistic:	7642.
Date:	Thu, 11 Oct 2018	Prob (F-statistic):	0.00
Time:	14:29:55	Log-Likelihood:	-80973.
No. Observations:	53576	AIC:	1.620e+05
Df Residuals:	53569	BIC:	1.620e+05
Df Model:	6		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	-3.6139	0.094	-38.551	0.000	-3.798	-3.430
average_rating	0.1185	0.022	5.427	0.000	0.076	0.161
review_delta	0.0416	0.004	9.747	0.000	0.033	0.050
images	0.1254	0.004	29.784	0.000	0.117	0.134
log_words	0.3031	0.005	59.478	0.000	0.293	0.313
log_times_reviewed	0.4092	0.003	145.136	0.000	0.404	0.415
log_user_followers	0.1319	0.002	70.752	0.000	0.128	0.136

Omnibus:	3352.761	Durbin-Watson:	2.005
Prob(Omnibus):	0.000	Jarque-Bera (JB):	4020.050
Skew:	0.649	Prob(JB):	0.00
Kurtosis:	3.340	Cond. No.	224.

R^2 : 0.467
MSE: 34

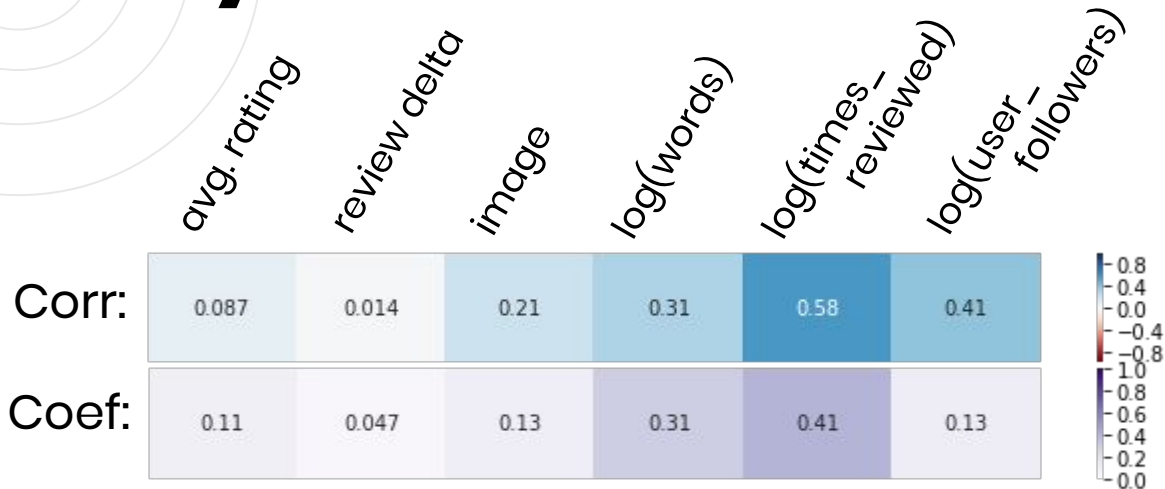


Model Improvement

- KFold
 - Robust across test/train splits
- Lasso & Ridge
 - Decrease in performance
- Polynomial Features
 - Unstable improvements
- Polynomial + ElasticNetCV
 - Instability + Worse Performance



Analysis: What Increase Likes?



- Read popular books
- Write longer reviews
- Use more images
- Get more friends



Next Steps

Natural language processing

