
CMPS 242 - Machine Learning (Final Report)

Analysis of Rating Dimensions Using Review Text

Bharath Nagesh
Madhu Shivashankaraiah
Numra Bathool
Tanay Parekhji

bnagesh@ucsc.edu
mshivash@ucsc.edu
nbathool@ucsc.edu
tparekhj@ucsc.edu

Git repository link :

https://github.com/madhusedu/topic_modeling.git

Abstract

Modern day systems take in both ratings and review texts from the user but generally consider only the ratings for performing different tasks, such as predictions and assessing the overall quality of a business. This according to us is not the optimal method. Ratings are highly subjective and what might be considered as favorable by one could be unfavorable for another. After a detailed observation of the review texts on IMDb, Zomato and Yelp we noticed that users implicitly provide information on their likes and dislikes. We plan to implement a system which takes in not only the rating but also the review text in order to form an opinion. We consider all the reviews of the users and then use data mining to summarize them. We also classify some reviews belonging to the same topic. form an opinion on the preferences of the users, then finally use those opinions in order to improve the overall rating quality for the business and perform better user predictions.

1. Problem Statement

Along with the expansion of the internet, a number of services such as e-Commerce and Review websites have emerged. A lot of people prefer using these services because of many factors, some of them being :

- **Reliability** : These services provide excellent delivery of product/information and then back it up with a good level of customer care. Also, the ratings provided on these sites aren't just a number that is

just floated to hoodwink the customer into buying the product/using the service making them very reliable

- **Ease of use** : All these services have a clean web/app interface which makes it very convenient for the people to use them.

The effect of the above factors can be seen in the growth of the industries, as indicated in Figure. 1 and Figure 2.

During holidays such as Thanksgiving the use of these services increase dramatically. The usage statistics is indicated in Figure 6.

The general method used to increase the reliability of a product/business/service is by collecting ratings reviews from the people who have used them. This facilitates future customers to form a better opinion of the product when deciding to whether or not use the product.

After a thorough observation of ratings and reviews provided on websites such as IMDb, Yelp, Amazon the following was observed. The overall rating for any product is formed as a weighted average of the ratings provided by all the users. This method completely disregards the review text provided by many of the users. The users generally provide important information which can be used to improve the overall rating of the product. For an example of this situation, we can consider reviews provided for a restaurant. An user may provide a rating of 3 stars for the restaurant and then provide a review which makes the restaurant actually worth 3.5 or maybe even 4 stars.

This problem that is presented above provided us with the opportunity to develop a system which uses Latent Dirichlet Allocation and Topic Modeling techniques in order to understand the reviews provided by the user. After doing so we attempt to develop an application which fine tunes the overall quality of ratings for any given product/business.

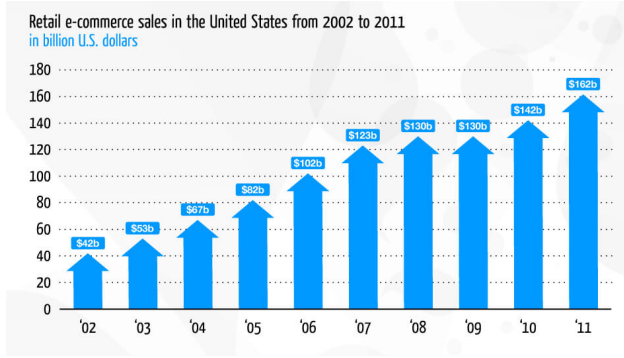


Figure 1. Retail sales in the USA between 2002 and 2011.

2. Algorithm Formulation

In order to perform all the operations we first started off by extracting the required text from the review data file. After this the extracted review text were processed into a format which can then be use within the code and identification numbers were assigned to every unique words. This list of unique words had to exclude words such as "a", "the" and so on. This was done to avoid wrong analyses of the extracted review text.

For our algorithm we decided to specify the number of topics that we would be using to group the set of documents into upfront (denoted by K). For our cases we defined 2 as the number of topics.

For the first part of algorithm we randomly assigned words to one the 2 topics that we have initialized. Along with generating this topic assignment list, we created two matrices which are :

- **Word-topic matrix :** in other words we created a count of the words being assigned to each of the 2 topics.
- **Document-topic matrix :** This is the matrix with the distribution of the topic assignment list.

In every document d , we observe every word w . We then chose a topic t with the $p(w | t) \times p(t | d)$, denoted by the following mathematical notations:

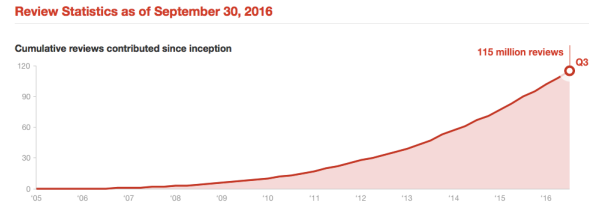


Figure 2. Yelp Review Statistics

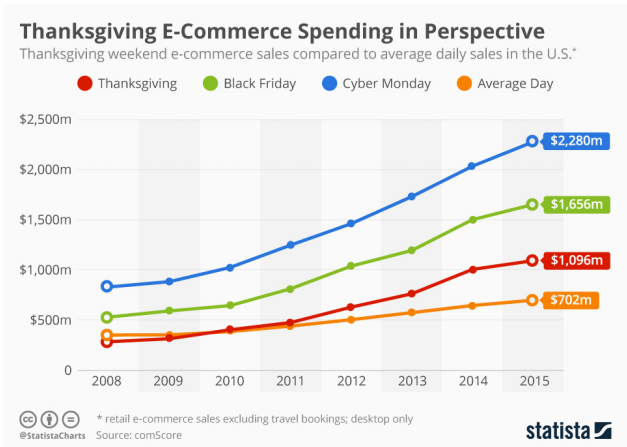


Figure 3. eCommerce sales during the Thanksgiving week.

$$p(z_i = j | z_{-i}, w_i, d_i) = \frac{C^{WT} w_i j + \eta}{\sum_{w=1}^W C^{WT} w j + W \eta} \times \frac{C^{DT} d_i j + \alpha}{\sum_{t=1}^T C^{DT} d_i t + T \alpha}$$

where,

$p(z_i = j)$: Probability that word i is assigned to topic j .

z_{-i} : Represents topic assignments of every other word.

w_i : Word ID of the i_{th} word.

d_i : Document containing the i_{th} word.

C^{WT} : Word-topic matrix.

$\sum_{w=1}^W C_{wj}^{WT}$: Total number of words in each topic.

C^{DT} : Document-topic matrix.

$\sum_{t=1}^T C_{d_i t}^{DT}$: Total number of words in document i .

η : Parameter to set the topic distribution for the words.

α : Parameter to set the topic distribution for the documents.

W : Total number of words in the set of documents.

T : Number of topics.

After drawing the new topic we update the topic assignment list with newly sampled topic for the word w

After that increment the word-topic and document-topic count matrices with the new sampled topic for word w .

After we complete learning the topics for r iterations iterations, we use the count matrices to obtain the word-topic distribution and document-topic distribution.

LDA assumes that each document is a mixture of all topics, thus after computing the probability that each document belongs to each topic we use this information to see which topic each of the documents belong to and possible words that are associated with each topic.

After understanding the computations we attempted to tune the parameters.

Since the starting point sampling point of dataset is chosen randomly, thus it makes sense to discard the first few iterations. Due to the fact that they most likely do not reflect the properties of distribution accurately. And another parameter is the number of iterations omitted during the training. This prevents the correlation between samples while iterating.

3. Feature Engineering

In order to perform the function of topic modeling we needed a dataset for which we used the review's dataset which was provided in the Yelp dataset. The dataset was found to be large in order to process and obtain results quickly. Due to time constraints it had to be sub-sampled to about 500 reviews. This gave the satisfaction that our results was obtained efficiently on a large enough dataset required to generalize our algorithm.

The primary feature that was required to perform the desired task was review text. Not much else was needed apart from this. Using the review text we created other parameters which we then used to create a topic model. Some of these were

- Word-topic matrix
- Document-topic matrix
- The various probabilities of occurrences

However we could not use the review text as obtained in the dataset. We had to process the extracted data so that the necessary functions could be applied to it in order to obtain the results.

In order to obtain and process the data we first had to convert the JSON file into a useable CSV file. Once this was done we sub-sampled the dataset into something that was smaller and would give us a flavor of what we could expect from a larger dataset.

Once the sub-sampling was complete we created a list of all the words within each review and then removed stop words i.e. words such as "a", "the" and so on... This helped us in forming a list with all the relevant and necessary words which could then be applied in our topic modeling algorithm.

After removing the stop words we went ahead and labelled the sentences i.e. we implemented topic modeling using the LDA algorithm and then classified each subtopic into groups such as Price, Time of the day, Service, etc...

We created specific methods within our code to implement the tasks.

4. Evaluation

Enter the evaluation and parameters used for the evaluation...

5. Results

[u'text'], [u'mr', u'haogie', u'institution', u'walking', u'seem', u'like',
u'haogie', u'voted', u'best', u'area', u'year', u'year', u'usually', u'border',
u'pod', u'superb', u'customer', u'service', u'miss', u'mario', u'machines',
u'that', u'fast', u'make', u'order', u'always', u'spot', u'fresh', u'veggies',
u'ords', u'ever', u'years', u'see'], [u'pros', u'italian', u'haogie', u'delicious',
u'failure', u'pizza', u'cheap', u'either', u'guess', u'home', u'says', u'ha',
u'union', u'rings', u'cost', u'25', u'74', u'm', u'thinking', u'reader', u'che',
u'sing', u'fairly', u'decent', u'haogies', u'italian', u'general', u'run', u'm',
u'second', u'steak', u'haogie', u'atrocious', u'cheese', u'add', u'cheese', u'1',
u'haogie', u'possibly', u'worst', u'valuable', u'haogie', u've', u'ever', u'eaten',
u'normally', u'reviews', u'establishment', u'unless', u'rating', u'exceptional',
u'trip', u'paid', u'24', u'59', u'two', u'whole', u'haogies', u'1', u'italic',
u'haogie', u'pay', u'1', u'add', u'cheese', u'haogie', u'traveling', u'eating',
u'money', u'declining', u'tomato', u'going', u'pay', u'dressing', u'side', u'p',
u'three', u'microscopically', u'thin', u'slices', u'meat', u'haogie', u'ham',
u'indicating', u'place', u'financial', u'trouble', u'blinding', u'view', u'or',
u'ing', u'ordered', u'without', u'lettuce', u'tomato', u'onions', u'dressing

Figure 4. Pre-Processed text that will be used for the implementation

[0], 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840,

Figure 5. Tokenized text

6. Challenges

During the course of the project we encountered a few significant challenges which we eventually were able to overcome. Some of the challenges that we encountered were.

- Understanding the LDA algorithm along with the parameters and the equations necessary for processing of the data
- Once we got a hold on the topic we had to arrive at a sub-sampled dataset which was large enough to get a good result

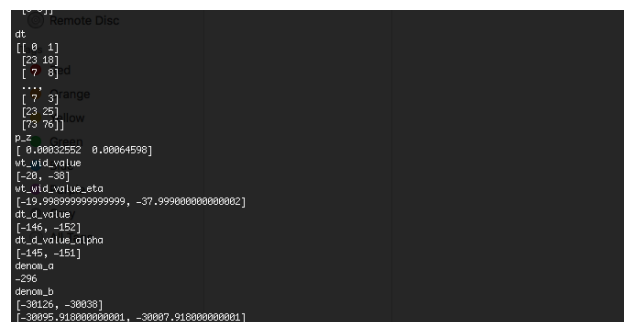


Figure 6. LDA Evaluation result

7. Learning objectives achieved

Over the course of the project over the last few weeks we achieved a good proficiency in many areas. The highlights of our learning include :

- We gained a good understanding of the Latent Dirichlet Allocation algorithm. Along with the algorithm itself we were able to understand it's functioning and the pro's and cons associated with it.

We were also able to get a good grip on it's implementation.

- It helped us discover and understand new packages available in python. Packages such as gensim and pyLDAvis.

8. Deliverables

We believe that 8 weeks was a short period for us to work on the project and that more could be done provided we had more time available. A few of the ideas that we had are :

- To create an API which could then be used to generate ratings based on the review text.
- To use the code against a larger dataset with more review texts.
- To use the code against different datasets. By this we mean that we intend to use the code against dataset obtained from other sites such as Amazon and eBay. This would help verify cross-platform functionality.

References

- [1] T. L. Griffith, M. Steyvers. *Finding Scientific Topics* PNAS, vol. 10, 6 April - 2004.
- [2] Ke Zhai, Jordan Boyd-Graber. *Online Latent Dirichlet Allocation with Infinite Vocabulary* In proceedings, *8th*

International Conference on Machine Learning, Atlanta, GA, 2013 JMLR : WCP volume 28. Copyright 2013 by the author

- [3] D.M. Blei, A. Y. Ng, M. I. Jordan. *Latent Dirichlet Allocation* Journal of Machine Learning Research. January, 2003
- [4] Blei, D. M., Griffiths, T. L., Jordan, M. I., Tenenbaum, J. B. *Hierarchical topic models and the nested Chinese restaurant process* In Advances in Neural Information Processing Systems 16. Cambridge, MA, USA: MIT Press. 2004
- [5] Erosheva, E. A. *Grade of membership and latent structure models with applications to disability survey data*. Unpublished doctoral dissertation, Department of Statistics, Carnegie Mellon University.
- [6] Hofmann, T. *Probabilistic Latent Semantic Analysis*. In Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence.
- [7] Zhuang.L, Jing.F, Zhu.X. *Movie Review Mining and Summarization*. CIKM'06, Arlington, VA, USA, 2006.
- [8] Hu.M, Liu.B *Mining and Summarizing Customer Reviews*. KDD '04, Seattle, WA, USA. 2004
- [9] McAuley.J, Leskovec.J *Hidden Factors and Hidden Topics: Understanding Rating Dimensions with Review Text* In Proceedings of the 7th ACM conference on Recommender systems, Hong Kong, China. 2013