

1 Solutions:

2

3 I have 3 features_matrices with 3 different group of features as explained below

4

5 1. Feature Engineering Approach 1:

6

7 a. Combined the user information of user_assessment_scores and user_interests taken an assessment.

8 b. On analysing the data, although all the given 10000 users have expressed interest, not all have taken the assessment. Hence there was loss of information in combining these two features.

9 c. Applied log transformation to the user_assessment_scores to achieve normal distribution of 'ratings' using unsupervised discretization with k-means algorithm. Used k=5.

10 d. Though the information was not present for all the users, the similarity was maintained.

11 e. For example, consider the users (user_handle: 12 and 4343), both the users have taken the assessment and scored close scores of 170 and 164 respectively.

12

13 2. Feature Engineering Approach 2:

14

15 a. Combined the user information of user_course_views and course_tags for each user.

16 b. Joining user_course_views with course_tags on course_id to get the course details for each user.

17 c. Creating a new Feature combining Course_tag and Level of the course. D = Degree, U = Undergraduate, P = Postgraduate.

18 d. Aggregating(sum) the number of seconds an user has viewed a course description.

19 e. Removing the entries where a user has viewed a course for '0' seconds.

20 f. Applying log transformation to the aggregated number of seconds to convert it into a rating.

21 g. Creating a new feature 'Rating':

22 0 - if the user has not viewed a course

23 1 - if the log_transformed_seconds in the range (0 - 1st quantile),

24 2 - if the log_transformed_seconds in the range (1st quantile - 2nd quantile)

25 3 - if the log_transformed_seconds in the range (2nd quantile - 3rd quantile)

26 4 - if the log_transformed_seconds in the range (3rd quantile - max)

27

28 h. Although there was little loss in information the similarity was considered.

29 for example, consider the top 3 similar users of user_handle : 10000 - user_10000

30 the course 'apex-absolute-beginner-guide-coding-salesforce' with 'beginner' tag.

31 has literally the viewed same courses.

32

33

34

35 3. Feature Engineering Approach 3:

36

37 a. The user_interest was the only feature with the information of all the given 10000 users. It was used to create a course at least once.

38 b. Used the interest_tag as the features for this model.

39

40

41 Solution for the questions given:

- 46
- 47 1. The metric chosen for similarity calculation is Cosine Similarity with Mea
- 48
- 49 a. Adjusted Cosine is invariant to scaling of the vector. Hence it is rob
- 50 while increasing the number of users or for more expressive features t
- 51 b. On calculating the coefficient based on all the elements in the vector
- 52 the similarity between users that may not have many features in common.
- 53 c. Hence it takes care of 2 main issues, dealing with uncommon features be
- 54 d. Adjusted Cosine is statistically similar to Pearson Correlation robust
- 55 e. Projection wise cosine similarity is dependent on only the direction o
- 56 f. The cosine similarities of a subset of the original data are the same
- 57
- 58 2. I would consider the following steps to scale the product.
- 59
- 60 a. Starting with storing the data, I have used SQLite for solving the pro
- 61 database or file storage like HDFS, Azure Data lake could be efficient.
- 62 b. Considering the data is stored on a distributed environment (Databrick
- 63 to implement parallel processing in the pipelines. For example consider d
- 64 be accomplished in parallel which is highly efficient.
- 65 c. Rendering the product to end-user should be a more containerized versi
- 66 or deploying the model on the cloud will provide substantial increase in
- 67 d. Using DAG engines like Apache Spark or Tensor Flow for numerical compu
- 68 process and computations internally through lazy evaluation. DAG engine w
- 69 graph assisting performance improvement. This could prove incredibly effi
- 70 which increases the computation at scale.
- 71 e. Perform computation on GPU
- 72 f. Perform random or stratified sampling (to ensure users from all groups
- 73
- 74 3. Considering the scale and diversity of the users the company serves, for a
- 75
- 76 a. Build a more robust API by changing the DB storage to distributed stor
- 77 b. As mentioned, definitely gather more data on the users (Job/student, p
- 78 Completed the Course/Not Completed/ Assessment Scores/ User Reviews/ e
- 79 expansive and expressive features. Also in the case of a Supervised Le
- 80 c. Track user journey through clickstream data to identify friction point
- 81 d. Develop end-to-end pipelines by covering more features like mentioned
- 82 e. Careful consideration of the points mentioned above for scaling the so
- 83
- 84
- 85
- 86
- 87
- 88