

RECOMMENDATION SYSTEM USING YELP DATASET

ABHAY KASTURIA
ADITYA PRIYADARSHI
GAUTAM VASHISHT
VARUN NANDU
XINGXING LIU

OUTLINE

- INTRODUCTION
- DATA PREPROCESSING
- METHODS
- EVALUATION AND RESULTS

INTRODUCTION

- RECOMMENDATION SYSTEM
- WHAT? WHY?

Recommendation Systems

ibuildings [i]
THE PHP PROFESSIONALS

amazon.com

Customers who bought this item also bought

The screenshot shows a section titled "Customers who bought this item also bought" on the Amazon website. It displays seven book covers with their titles and prices:

- Test Driven Development: By Example by Kent Beck \$49.98
- Design Patterns: Elements of Reusable Object-Oriented Software by Erich Gamma \$38.70
- Refactoring to Patterns by Joshua Kerievsky \$41.93
- Working Effectively with Legacy Code by Michael C. Feathers \$38.70
- Patterns of Enterprise Application Architecture by Martin Fowler \$37.44
- The Pragmatic Programmer: From Journeyman to Master by Andrew Hunt \$32.25
- Code Complete: A Practical Handbook of Software Construction by Steve McConnell \$31.49

last.fm
the social music revolution

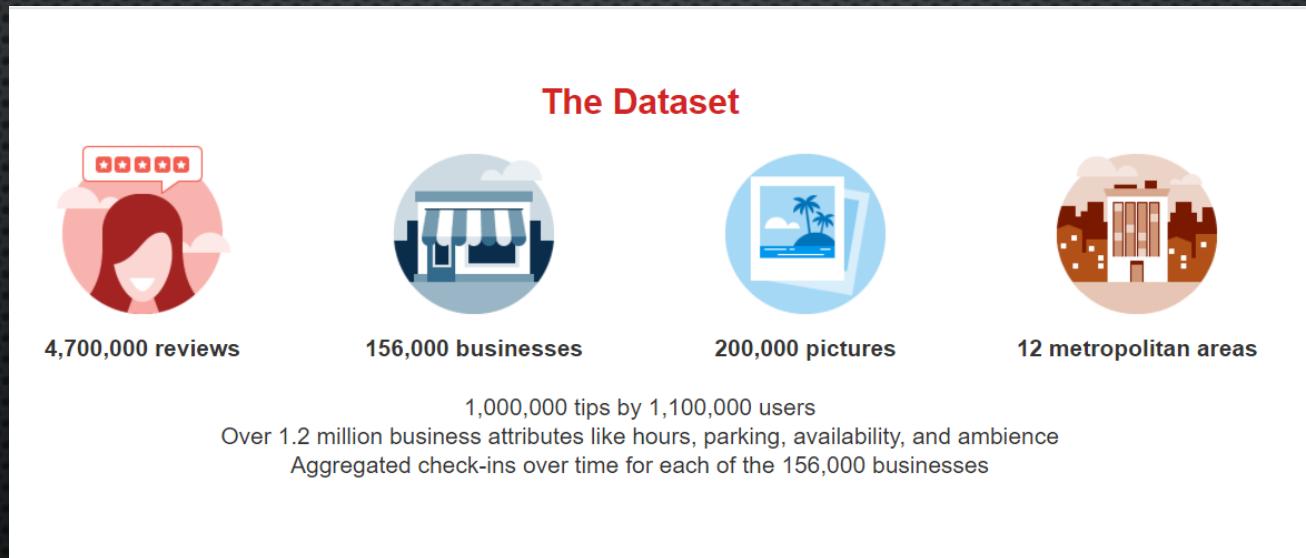
Music Recommended by Last.fm

The screenshot shows a section titled "Music Recommended by Last.fm" on the Last.fm website. It displays four recommendations with small images and names:

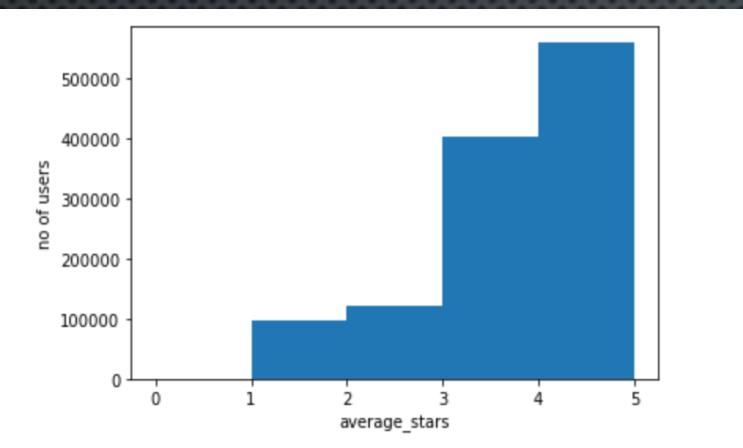
- John Patitucci**
Similar to: Jaco Pastorius, Victor Bailey, Chick Corea
- Yellowjackets**
Similar to: Chick Corea, Weather Report, Dave Weckl
- Stanley Clarke**
Similar to: Jaco Pastorius, Marcus Miller, Weather Report
- Weather Report**
Similar to: Jaco Pastorius, Chick Corea, Marcus Miller

PROJECT

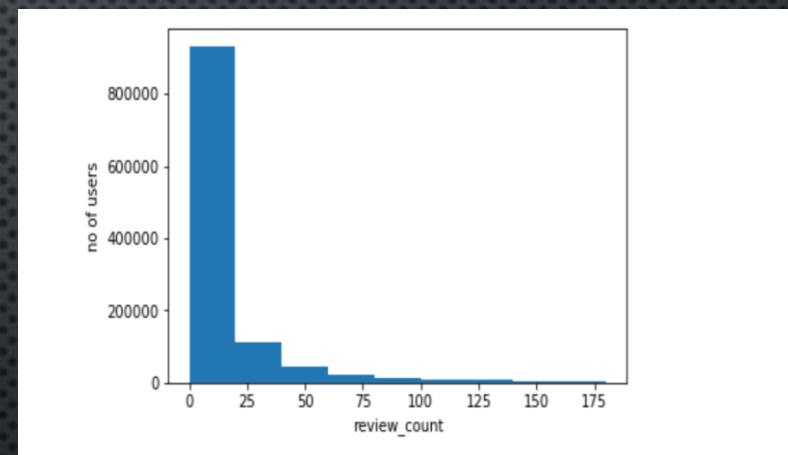
- USE YELP DATASET TO RECOMMEND RELEVANT BUSINESSES TO USERS
- BASED ON THE ACTIVITY OF USERS ON YELP
- ACTIVTY = REVIEWS GIVEN BY USER TO BUSINESS
- EVALUATE SEVERAL DATA MINING METHODOLOGIES USED FOR RECOMMENDATIONS



BASIC DATA ANALYSIS

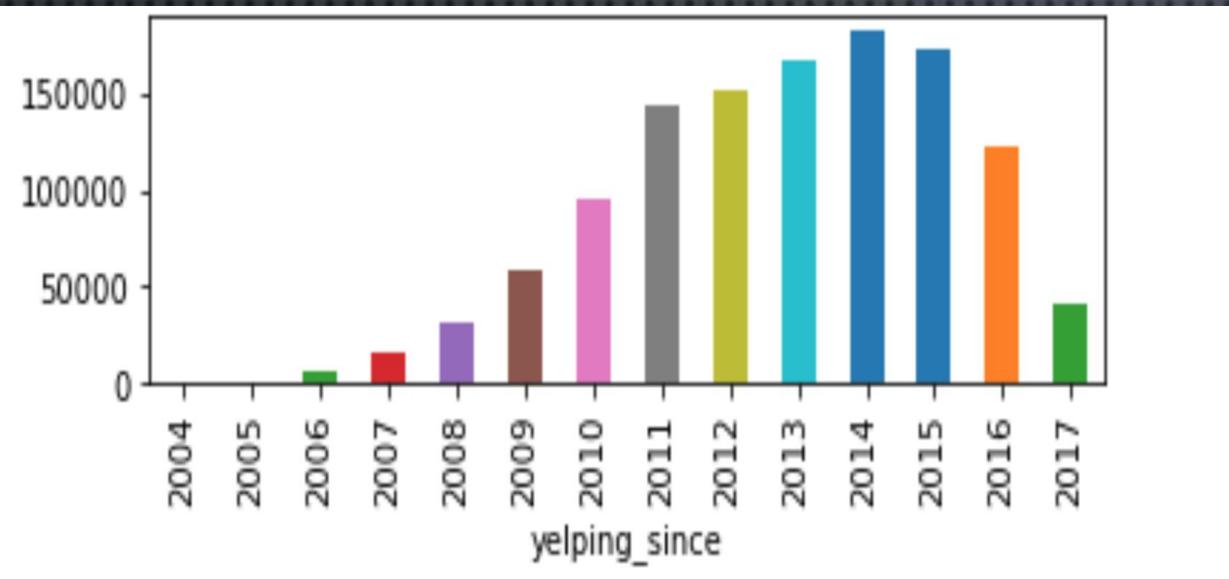


Average Star Ratings

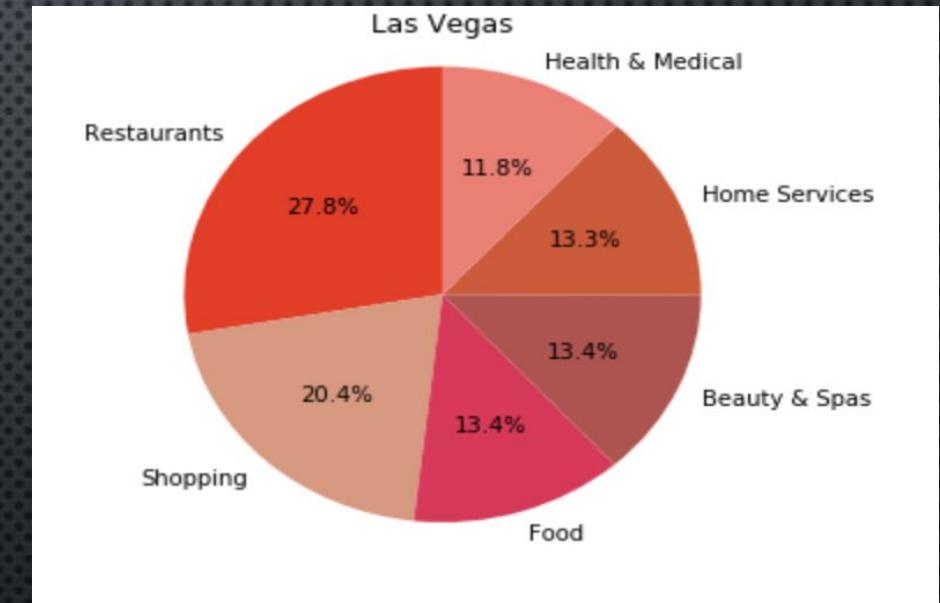


Average #Reviews

BASIC DATA ANALYSIS



User Growth



Popular Categories in Las Vegas

DATA PREPROCESSING

- DATA FILTERING:
 - USERS WITH REVIEWS ≥ 20 ~20K RECORDS
 - BUSINESS WITH CATEGORY = RESTAURANT ~ 48K RECORDS
- DATA TRANSFORMATION:
 - USERID AND BUSINESSID MAPPED TO SERIAL NUMBERS
- DATA DIVISION:
 - BASED ON TIME
 - TRAINING DATA – REVIEWS TILL 2015 ~ 900K
 - TESTING DATA – REVIEWS AFTER 2015 ~ 19K
 - TESTING DATA REVIEW CATEGORIZATION - POSITIVE(>3.7), NEGATIVE(<2.6), NEUTRAL(REST)

METHODS

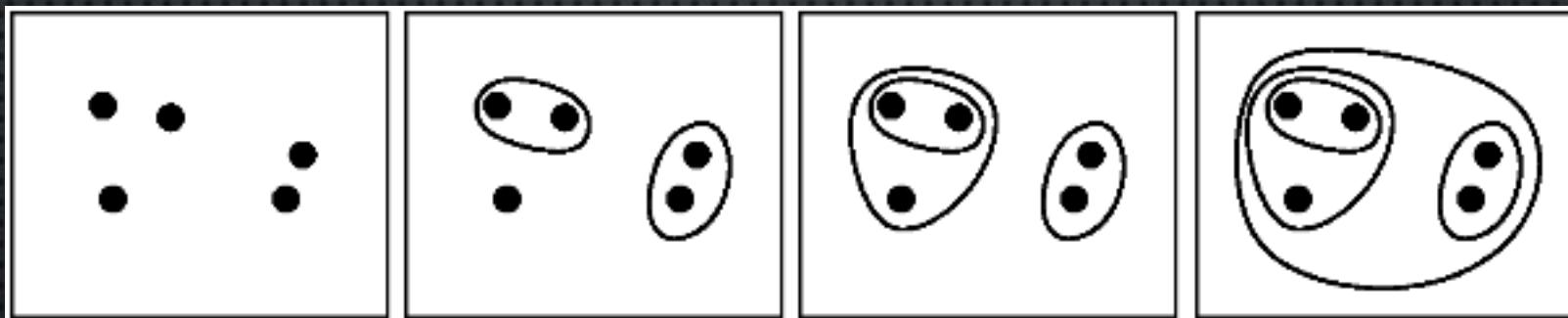
- CLUSTERING
 - AGGLOMERATIVE
- COLLABORATIVE FILTERING
 - USER – USER SIMILARITY BASED
 - MATRIX FACTORIZATION
- DEEP LEARNING

CLUSTERING - APPROACH

- ASSUMPTIONS – USERID, LOCATION AND TYPE OF BUSINESS
- USER PREFERENCE VECTOR OF TOP 15 CATEGORIES BASED ON RATING – WHAT IF NO RATING?
- BUSINESS VECTOR – MUST CONTAIN 3 OF TOP 15 CATEGORIES.
- USE AGGLOMERATIVE HIERARCHICAL CLUSTERING TO CLUSTER USERS

AGGLOMERATIVE HIERARCHICAL CLUSTERING

Number Of Clusters?



RECOMMEND BUSINESSES

- CALCULATE AVERAGE PREFERENCE VECTOR
- CALCULATE DOT PRODUCT WITH EACH BUSINESS VECTOR
- SORT THE BUSINESSES BASED ON DECREASING VALUE OF THE ABOVE DOT PRODUCT
- TAKE OUT FIRST 20 BUSINESSES AS YOUR TOP RECOMMENDATIONS

CLUSTERING - RESULTS

- HIT RATIO
 - POSITIVE(>3.7), NEGATIVE(<2.6), NEUTRAL(REST)
- NORMALIZED DISCOUNTED CUMULATIVE GAIN (NDCG)

City	Hit Ratio	NDCG
Las Vegas	0.58	0.78
Toronto	0.61	0.75
Phoenix	0.62	0.80
Charlotte	0.55	0.77
Edinburgh	0.60	0.81
Pittsburgh	0.44	0.78
Montral	0.57	0.80
Scottsdale	0.69	0.91
Cleveland	0.64	0.89
Madison	0.64	0.78
Stuttgart	0.5	0.71
Mesa	0.41	0.52
Tempe	0.35	0.69
Henderson	0.65	0.88
Mississauga	0.22	0.63
Average	0.55	0.78

USER-USER BASED COLLABORATIVE FILTERING

- CREATED UTILITY MATRIX FOR USERS VS BUSINESS FILLED WITH RATING
- CENTRALIZED EACH RATING VECTOR OF USER TO MEAN ZERO
- REPLACE THE MISSING VALUES WITH ZERO
- AVERAGED RATING VALUE WOULD BE 0, POSITIVE OR NEGATIVE WHERE 0 REPRESENTS MEAN RATING OR MISSING RATINGS AND POSITIVE REPRESENTS HIGHER RATINGS THAN MEAN RATING

RATING PREDICTION

$$predr_{(i,k)} = \frac{\sum sim(u_i, u_j) * r_{(j,k)}}{\#users(j)}$$

- PREDICTION FOR USER I FOR BUSINESS K IS GIVEN BY THE FORMULA ABOVE.
- PREDICTED THE RATINGS THAT USER WILL GIVE TO THE NEW BUSINESS
- FIND TOP K SIMILAR USERS WHICH HAVE ALREADY RATED THE ITEM ‘I’ AND RETURNED THE WEIGHTED AVERAGED OF RATINGS WITH WEIGHT AS SIMILARITY

RECOMMENDING TOP 10 BUSINESS TO USER

- PREDICT RATING FOR BUSINESS FOR THE USER WHICH ONE HAVE NOT USED YET
- SORT THE ARRAY ON BASIS OF RATING
- RETURN THE TOP 10 BUSINESS

COLLABORATIVE FILTERING MATRIX FACTORIZATION USING ALS

$$pred_{ui} = q_i^T \cdot p_u$$

- Q – VECTOR OF BUSINESS AND P- VECTOR OF USERS
- PREDICTED THE RATINGS BASED ON THE DOT PRODUCT OF THE USER AND THE BUSINESS VECTOR

$$\min_{q,p} \sum_{(u,i)} (r_{ui} - q_i^T p_u)^2 + \lambda(\| q_i \|_2^2 + \| p_u \|_2^2)$$

- BOTH Q AND P ARE UNKNOWNS, EQUATION IS NOT CONVEX.
- ALS FIXES ONE AND OPTIMIZES THE OTHER AND ROTATE BETWEEN FIXING THE Q AND FIXING THE P. THUS, THE NAME – ALTERNATING LEAST SQUARES.
- THIS ENSURES THAT EACH STEP DECREASES EQUATION UNTIL CONVERGENCE.

COLLABORATIVE FILTERING MATRIX FACTORIZATION USING ALS

$$\min_{q,p} \sum_{(u,i)} (r_{ui} - q_i^T p_u)^2 + \lambda(\| q_i \|_2^2 + \| p_u \|_2^2)$$

- PARAMETERS FOR ALS:
 - LAMBDA - REGULARIZATION PARAMETER
 - ITERATIONS – MAX IF CONVERGENCE NOT REACHED
 - RANK - THE NUMBER OF FEATURES TO USE (AKA THE NUMBER OF LATENT FACTORS).
- PARAMETERS FOR OUR ALS MODEL:
 - LAMBDA = 0.15
 - ITERATIONS = 20 (SLOW COMPUTERS)
 - RANK = 100 (LOW MEMORY COMPUTERS)

COLLABORATIVE FILTERING MATRIX FACTORIZATION USING ALS

$$\min_{q,p} \sum_{(u,i)} (r_{ui} - q_i^T p_u)^2 + \lambda(\| q_i \|_2^2 + \| p_u \|_2^2)$$

- RESULTS:

	MSE - Training	MSE - Testing
Non-Reg	0.02	2.2
Reg	0.04	1.4

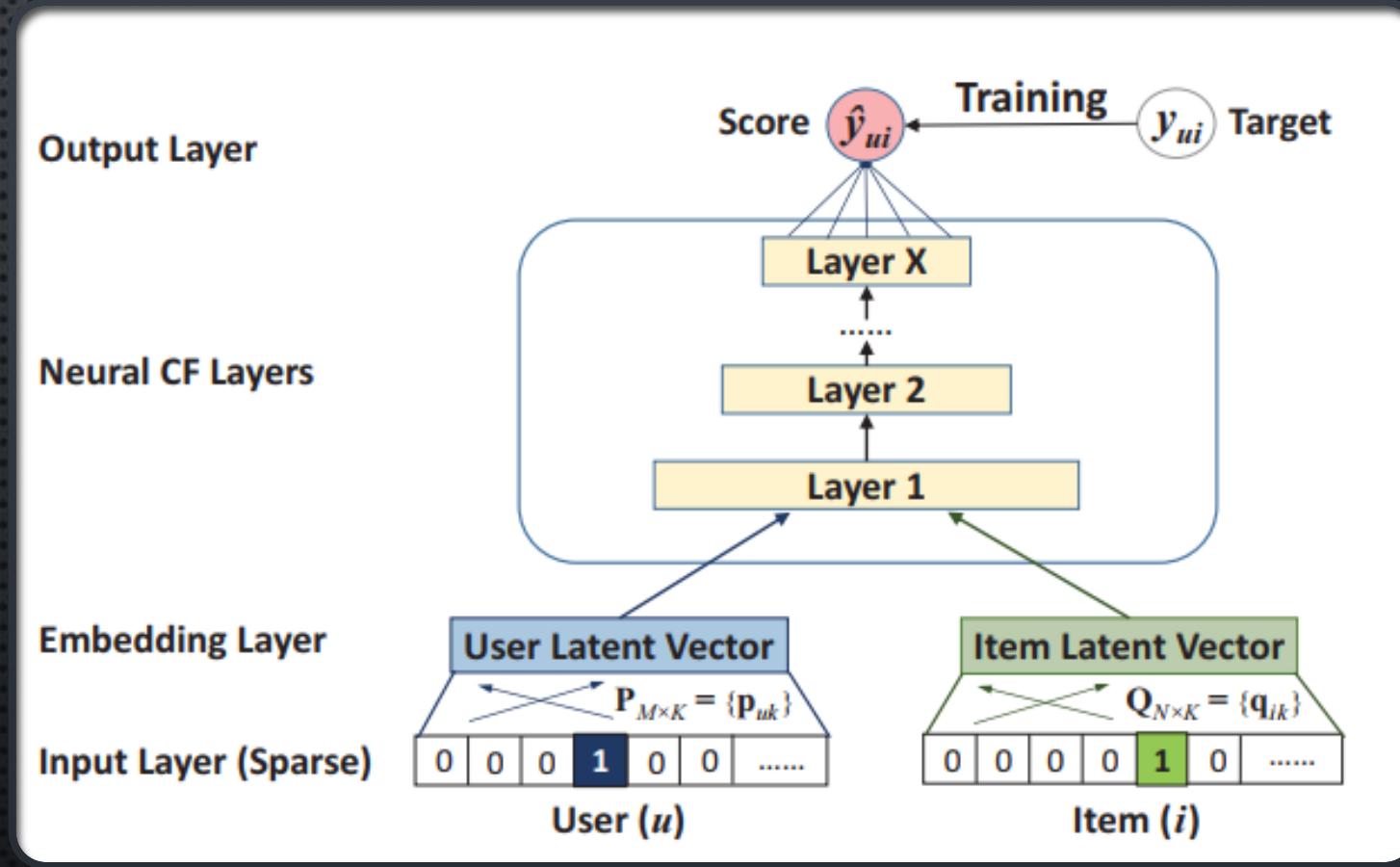
	Hit Ratio(Threshold = 1.5)
Reg	0.77

DEEP LEARNING - APPROACH

- USES MULTI-LAYER PERCEPTRON (MLP) TO LEARN USER-ITEM INTERACTION
- USES IMPLICIT FEEDBACK RATHER THAN EXPLICIT REVIEWS

$$y_{ui} = \begin{cases} 1, & \text{if interaction (user } u, \text{ item } i) \text{ is observed} \\ 0, & \text{otherwise} \end{cases}$$

DEEP LEARNING - ARCHITECTURE



DEEP LEARNING - RESULTS

- TRAINED MODEL USING 4 HIDDEN LAYERS WITH [64,32,16,8] FACTORS.
- USED HIT RATIO AND NDCG TO MEASURE QUALITY
- BEST MODEL – HIT RATIO (0.8197) AND NDCG (0.6220)

CONCLUSION

	Hit Ratio
CF - ALS	0.77
DNN	0.81
Clustering(avg across cities)	0.55

- WE SEE THAT DNN GIVES US THE BEST HIT RATIO

FUTURE WORKS

- UTILIZE TEXTUAL AND IMAGE BASED DATA TO HELP IMPROVE THE RECOMMENDATIONS.
- UTILIZE HYBRID RECOMMENDATION SYSTEM TECHNIQUES TO USE USER BEHAVIOUR AND USER CHARACTERISTICS(DEMOGRAPHICS OF USERS) TO IMPROVE RECOMMENDATIONS.