

Index:
From Before:
Text Preprocessing

Statistical Properties and Evaluation

Boolean Retrieval

Vector Space Model

Relevance Feedback

Statistical Language Modeling

Document Prior

Page Rank

New:

Page Rank

Recommendation Systems

Text Classification

Text Clustering

Page Rank

Recommendation Systems

Algorithms

- Collaborative Filtering
- Content-Based
- Hybrid

Collaborative Filtering

- Like users rate similarly
- k-nearest
 - choose the k most similar

* Evaluating Predictions: RMSE

$$\sqrt{\frac{1}{S} \sum_{(u,i) \in \text{test}} (p_{u,i} - r_{u,i})^2}$$

no under/over

P ratings predicted
r true rating
(u,i) < test missing ratings
S total predicted

- Neighborhood Selection

Cosine Sim

$$w_{u,u} = \frac{\sum_{i=1}^n r_{u,i} \times r_{u,i}}{\sqrt{\sum_{i=1}^n r_{u,i}^2} \times \sqrt{\sum_{i=1}^n r_{u,i}^2}}$$

Rating Prediction

$$p_{u,i} = \frac{\sum_{u=1}^k w_{u,u} r_{u,i}}{\sum_{u=1}^k w_{u,u}}$$

- Pearson Correlation

- model the deviation

$$w_{u,u} = \frac{\sum_{i=1}^n (r_{u,i} - \bar{r}_u)(r_{u,i} - \bar{r}_u)}{\sqrt{\sum_{i=1}^n (r_{u,i} - \bar{r}_u)^2} \times \sqrt{\sum_{i=1}^n (r_{u,i} - \bar{r}_u)^2}}$$

$$p_{u,i} = \bar{r}_u + \frac{\sum_{u=1}^k w_{u,u} (r_{u,i} - \bar{r}_u)}{\sum_{u=1}^k |w_{u,u}|}$$

r
Average

sub ratings:

- 1 - perfectly correlated
- 0 - not
- 1 - inverse

Optimizations Improving Predictions

- Penalize universally liked movies

$$IUF(j) = \log \frac{m}{m_j} \quad \text{total A users} \quad \text{users rated item j}$$

- multiply original ratings by IUF during weight calc

- Case Amplification

$$w_{u,u} = w_{u,u} \cdot |w_{u,u}|^p, \quad p = 2.5, \geq 1$$

- favors high weights, punishes low weights

Euclidean dist:

$$d = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}, \quad w = 1/d+1$$

Jaccard Sim

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \rightarrow \text{binary ratings}$$

User Based - Summary

- 1) find k nearest
- 2) use ratings from k-nearest to predict for the active user

Item-Based Collaborative Filtering

- Similar Items rated similarly - reverse users & items in previous

Cos sim

$$w = \sum_{i=1}^n r_{u,i} \times r_{u,i}$$

Pearson: $\sum (r_{u,i} - \bar{r}_i)(r_{u,j} - \bar{r}_j)$

$$\sqrt{\sum_{j=1}^n r_{u,j}^2} \times \sqrt{\sum_{j=1}^n r_{u,j}^2}$$

$$\sqrt{\sum_{j=1}^n (r_{u,j} - \bar{r}_i)^2} \sqrt{\sum_{j=1}^n (r_{u,j} - \bar{r}_u)^2}$$

Adj's cos sim

Replace \bar{r}_i w/ \bar{r}_u

User vs Item

- Item more efficient
- Item more stable - users have multiple interests
- Item based - little diversity or surprise - obvious and boring

Cold-Start Problem

- new item/users have no historical ratings

Content-Based

content compare to user profile

IR applies - treat user as a query, item as document

★ K-Nearest - Probability improvements

- center your data
- more points less noisy, orders of magnitude for ratings
- Smoothing!

$$\text{Linear} - r_u = \alpha \cdot r_u + (1 - \alpha) \cdot g \Rightarrow g \text{ global mean} \quad r_u \text{ is } r_u \text{ average rating}$$

$$r_u = \frac{\sum_{j=1}^n r_{u,j}}{n_u} \Rightarrow \text{K nearest neighbors to } g$$

$$\frac{1}{n_u} \cdot g \Rightarrow \text{Dirichlet} \quad n_u \neq \text{ratings of user } u$$

- got a shrunken mean

Text Classification



- Assign Categories to text documents
- Vector Space - word = component

Like Proccio

$$\frac{1}{|D_c|} \sum_{d \in D_c} V(d)$$

D_c set of all documents in class c , $V(d)$ vector of d



K-nearest

- keep all training docs
- k docs that are most similar to the new doc
- assign category that is most common amongst neighbor docs

similarity functions

- euclidean, KL divergence, Dot product, cosine similarities

$$\hat{y}(x) = \frac{1}{\sum_{i=1}^k \text{sim}(x, x_i)} \sum_{i=1}^k \text{sim}(x, x_i) y_i$$

$$\begin{cases} 1 \text{ if } \hat{y}(x) > 0.5 \\ 0 \text{ otherwise} \end{cases}$$

★ Naive Bayes Text Classification

- query likelihood, doc prior

c - doc of all docs in

$$P(C, d) = P(C|d)P(d) = P(d|C)P(C) \text{ prior}$$

$$P(C|d) = \frac{P(d|C)P(C)}{P(d)}$$

one class
d - query
- compare likelihood

$$C = \operatorname{argmax}_j P(C_j) \prod_{i=1}^n P(t_i|C_j)$$

$$P(t_i|C_j) = \frac{t(t_i, C_j)}{|C_j|} P(C_j) = \frac{\text{does in } C_j}{\text{does}} \xrightarrow{\text{smoothing}} P(t_i|C_j) = \frac{t(t_i, C_j) + 1}{|C_j| + |V|}$$

add-one, etc

	in the class	not in the class
predicted to be in the class	true positives (TP)	false positives (FP)
predicted to not be in the class	false negatives (FN)	true negatives (TN)

$$\text{Precision} = TP / (TP + FP)$$

$$\text{Recall} = TP / (TP + FN)$$

$$F = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall})$$

$$\text{Accuracy} = (TP + TN) / (TP + FP + FN + TN)$$

Text Clustering

unsupervised learning
inferred from data

vs classification
supervised
user defined

* K-means $\vec{\mu}(c) = \frac{1}{|c|} \sum_{\vec{x} \in c} \vec{x}$

- iterative - based on centroids
- reassignment based on distance to current cluster centroids
- recompute centroids based on new membership

- guaranteed to converge
- not optimal - need good seeds
- only on linear data

$$k \approx \sqrt{n/2}, n \text{ is data points}$$

Find Sample

1) Summarize algorithms

a) User-based. Like users rate similarly

b) Item-based. Similar items are rated similarly

2) Major disadvantage to RMSE?

| Subjective Metric \neq RMSE

- doesn't correlate w/ user-satisfaction

Can't differentiate between above & below ground truth
- no under/over

3) K-means doesn't work where?

Non-gaussian data - ex: concentric circles.

The centroids are the same, but the algorithm won't find that

vs

4) $C = \text{America}$, $\bar{C} = \text{everything else}$.

5 training, 1 test: 6 total docs. Add one smoothing

$$\frac{P(d|C) P(C)}{P(d)}$$

$$P(\text{America}) = 5/12 \leftarrow 11 \text{ total words} + 1 = 5/12$$

$$P(\text{Calif}) = 1/12$$

$$P(\text{Wash}) = 1/12$$

$$P(\text{Oregon}) = 1/12$$

$$P(\text{Tokyo}) = 2/12$$

$$P(\text{Japan}) = 1/12$$

$$P(C) = \left(\frac{5}{12}\right)^3 \cdot \left(\frac{1}{12}\right) \cdot \left(\frac{1}{12}\right)^{3/4}$$

$$P(C) = 3/4 \leftarrow \text{doc prior} - 3/4$$

$$P(\bar{C}) = 1/4 \leftarrow 1 \text{ doc not in } C$$

$$P(d|C) = \prod_t P(t|C), t \in d[acc]$$

$$= P(\text{America}|C)^3 \cdot P(\text{Tokyo}|C) \cdot P(\text{Japan}|C)$$

$$\Rightarrow |C_j| = 7, |V| = 6$$

$$= \left[\frac{5}{13} \right]^3 \left[\frac{1}{13} \right] \left[\frac{1}{13} \right] \left[\frac{3}{4} \right] = 2.52 e^{-4}$$

$$P(d|\bar{C}) = \left[\frac{2}{7} \right]^3 \left[\frac{3}{7} \right] \left[\frac{3}{7} \right] \left[\frac{1}{4} \right] = 7.13 e^{-4}$$

$$C_j = 4, |V| = 3$$

$$P(d|\bar{C}) > P(d|C) \Rightarrow \text{in } \bar{C}$$

$$P(t_i|C_j) = \frac{f(t_i, C_j) + 1}{|C_j| + |V|}$$