

Clustering

1. **Each dimension** corresponds to a **feature/attribute**

2. **Why Clustering Hard** (1) Number of clusters is unknown

Arbitrary shapes and sizes (3) Quality of clustering result: a. depends on the distance measures b. measured by the ability to find hidden patterns (4) applications in high dimensional space 3.

Application: Sloan Digital Sky Survey & Music CDs & Cluster Documents (Find topics)

4. **Distance Measure:** Non-negativity: $d(x,y) \geq 0$, Identity: $d(x,y)=0$ if and only if $x=y$, symmetry, $d(x,y) = d(y,x)$,

triangle inequality: $d(x,y) \leq d(x,z) + d(z,y)$ 5. **Examples: Lp-norm:**

$$d(x,y) = \left(\sum_i (x_i - y_i)^p \right)^{\frac{1}{p}} \quad \text{L2-norm=Euclidean Distance}$$
$$d(x,y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

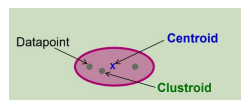
L1-norm: Manhattan Distance $d(x,y) = |x_1 - y_1| + |x_2 - y_2| + \dots +$

Cosine: $\text{cosine}(x,y) = \frac{\sum_i x_i y_i}{\sqrt{\sum_i x_i^2} \sqrt{\sum_i y_i^2}}$ **Jaccard:** $J(A,B) = \frac{|A \cap B|}{|A \cup B|}$ **Edit**

Distance: distance that convert one string to another **Hierarchical**

Clustering: repeatedly combine two nearest clusters **How to**

represent Cluster? (1) Euclidean space: centroid = average of its data points (2) Non-Euclidean Space: Clustroid= point "closest" (smallest maximum/average/sum of squares distance) to other points



$$\min_c \sum_{x \in C} d(x,c)^2$$

Intercluster distance = minimum of

distances between any two points, one from each clusters **Diameter:**

The maximum distance of two points in the cluster **Aver distance**

between points in the cluster **Density-based:** Diameter or Average

distance / the number of points in the cluster **When to Terminate** (1)

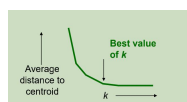
Pre-determined number of clusters (2) merging two clusters -> 'bad' cluster (diameter of merged cluster > threshold, diameter > average diameter by a wide margin / Density of clusters < threshold)

Complexity of hierarchical clustering: Naive: $O(n^3)$ /Priority

Queue: $O(n^2 \log n)$ **K-means Algorithm:** Process: 1. place points to nearest clusters (centroid) 2. update centroid of k clusters 3.

reassign points to closest centroid 4. repeat 2 and 3 until convergence

Select K:



Complexity-K-means:

- Complexity is $O(n \cdot k \cdot m)$
- n = number of points
- k = number of clusters
- m = number of iterations

BFR Algorithm: Variant of k-means to handle very large data sets in

high dimensions./ points in clusters are normally distributed around a

centroid in a **Euclidean space Memory** $O(\text{clusters})$ 1. Points are read

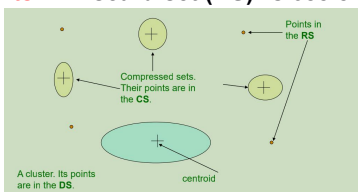
from disk to memory in chunk 2. Process: (1) Select initial k centroids

from the first chunk (Take k random points / Take a small random

sample and cluster optimally / Pick a random point, and then k-1

more points (each as far from the previous points as possible) 3.

Points: 1.Discard set (DS): Close enough to a centroid are



summarized 2.

Compression set (CS):

Summarized, but not

assigned to a cluster

3.**Retained set (RS):**

Isolated points **Summarizing Points** N (number of points) **SUM**(sum of points in ith dimension) **SUMSQ**(sum of squares of coordinates in ith dimension) **2d(number of dimensions)+1** -> any size cluster

Centroid (Average in each dimension) SUM_i/N **Variance** $(\text{SUMSQ}_i/N) - (\text{SUM}_i/N)^2$

Mahalanobis

Distance

For a point (x_1, \dots, x_n) and centroid (c_1, \dots, c_n)

$$d(x,c) = \sqrt{\sum_{i=1}^n \left(\frac{x_i - c_i}{\sigma_i} \right)^2}$$

Normalized Euclidean from Centroid

1. Compute MD (point & cluster centroid) 2. Choose cluster with least MD 3. Add point to cluster (MD < threshold)

Combine 2 CS into one: 1. Compute the variance of combined cluster 2. N , SUM and SUMSQ make calculation quickly 3. Combine if the combined variance < threshold

CURE Algorithm 1. Difference: (1) BFR and k-means assum

cluster are normally distributed in each dimension [Axes are fixed, Ellipses at an angle are not OK] (2) CURE assumes Euclidean

distance & allow clusters of any shape & Use a collection of

representative points to represent clusters] 2. **Process:** (1) **Pass 1:** a.

Pick a random sample of points that fit in main memory (2) Cluster

these points hierarchically(group nearest points/clusters) (3) Pick

representative points: for each cluster, pick a sample of points, as

dispersed as possible / For the sample, pick representatives by

moving them 20% toward the centroid of the cluster (2) **Pass 2:** 1.

Rescan the whole dataset and visit each point p in the dataset 2.

Place it in the "closest cluster"

Recommender System

1. **Long tail:** a near-limitless selection 2. **Types of**

Recommendations: a. Editorial (list of favorites / essential items)

b. Simple aggregates (Top10, Most popular, Recent Uploads)

c. Personalized 3. **Utility Function:** $C \times S \rightarrow R$ [C = sets of customers/

S = set of items / R = set of ratings] 4. **Key Challenges:** a. Gathering

"known" ratings for matrix b. Extrapolating unknown ratings from

known ones c. Evaluating extrapolation methods (how to measure

the (success/performance) of recommendation methods) 5.

Extrapolating Utilities: a. Utility matrix U is sparse [no ratings / cold

start] b. Recommender system seen as a function (1) Given -> User

model (i.e. Ratings, preferences, demographics 人口统计) & Items

(with or without description of item characteristics) (2) Find

Relevance Score **Content-Based Recommender Systems 1.Main**

idea: recommend to customer(C) **items** that are similar to previous

items rated highly by C 2. **Examples:** Movie recommendations / web,

blogs, news 3. **Item Profiles:** Profile is a set of features (value->

boolean or real-valued), use TF-IDF to pick important features

4. **User Profiles:** a. vector to describe user preferences b. Predict

preference of x for items that has not rated

Given user profile x and item profile i , estimate

$$\text{utility}(x,i) = \cos(x,i) = \frac{x \cdot i}{\|x\| \cdot \|i\|}$$

5. **Pros:** 1. No need for data

on other users 2. Able to

recommend to users with unique tastes 3. Able to recommend new &

unpopular items 4. Able to provide explanations 6. **Cons:** 1. Finding

the appropriate features is hard 2. Cold start problem for new users

3. Overspecialization [Never recommend items outside user's

content profile / People might have multiple interests] 4. Unable to

exploit quality judgements of other users. **Collaborative Filtering**

1.Main idea: Estimate x's ratings based on the ratings of users in N whose ratings are 'similar' to x's ratings

2. Similar Users: a. Jaccard similarity measure b. Cosine similarity c. Pearson correlation coefficient

2. Rating Predictions

$$\text{sim}(x, y) = \cos(r_x, r_y) = \frac{r_x \cdot r_y}{\|r_x\| \cdot \|r_y\|}$$

▪ S_{xy} = items rated by both users x and y

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_x} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_y} (r_{ys} - \bar{r}_y)^2}}$$

▪ Use average of their ratings $r_{xi} = \frac{1}{k} \sum_{y \in N} r_{yi}$

weighted measures $r_{xi} = \frac{\sum_{y \in N} s_{xy} \cdot r_{yi}}{\sum_{y \in N} s_{xy}} \quad s_{xy} = \text{sim}(x, y)$

3.Collaborative Filtering: (1) Pros: Work for any kind of item (no feature selection needed) **(2) Cons:** 1. Code Start (Need enough users in the system to find a match) 2. Sparsity (Hard to find users that have rated the same items) 3. First rater (Cannot recommend an unrated item) 4. Popularity bias (tend to recommend popular item)

4.Hybrid Recommenders: Add content-based method to CF

▪ **Root-mean-square error (RMSE)**

$$\sqrt{\sum_{xi} (r_{xi} - r_{xi}^*)^2}$$

r_{xi}^* is predicted
 r_{xi} is the true rating of x on i

5. Evaluation Predictions:

6.complexity: Find k most similar customers $\rightarrow O(k \cdot |C|)$ **Solutions:** 1. Near-neighbor search in high dimensions (LSH) 2. Clustering 3.

Dimensionality reduction **7.Tips: Add Data** (Leverage all data)

Analyzing Large Graphs

1.Community Detection: a.inter-group-sparse b.intra-group-dense

2.Edge Betweenness: a.how important an edge b.number of shortest path through an edge c.inter-community edge has higher betweenness

3.Girvan-Newman Algorithm: (1) Breadth first search of graph start from one node (2) Top down: root A=1, next label is sum of parent. (3) Bottom up: [Credit (Z) * label (Pi) divided by sum of labels of P1...,Pk] [1 + credit on edge below it] (4) Repeat the calculation for every node

$$B(a, b) = \sum_{x, y \in V} \frac{|SP(x, y) \text{ that include } (a, b)|}{|SP(x, y)|}$$

as the root and sum contribution (5) Divide by 2 to get final edge betweenness

4.Scalability issues: m edges and n nodes \rightarrow Final betweenness score for every edge $O(mn)$ & Recalculation takes $O(m^2n)$

5.Scaling up G-N Algorithm: 1.Determine which edges need recalculate betweenness when an edge is removed. (When e is removed, betweenness (e') needs to be recalculated only if e' is in the same connected component as e 2.Not much pruning if a component is large. **6.Link Analysis:** (1) Types of links: Referential - Click here and get back home (2) Informational-Click here to get more detail **7.PageRank:** Links as votes (A page is more important if it has more links) **8.Matrix Formulation:** (1) M (2) Rank vector $r \rightarrow r$ is the important score of page i. (3) Flow equation: $r = M \cdot r$

9. Power iteration Method:

- Power iteration – simple iterative scheme
 - Suppose there are N web pages
 - Initialize $r^{(0)} = [1/N, \dots, 1/N]^T$ & compute $r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$ d_i out-degree of node i
 - Iterate: $r^{(t+1)} = M \cdot r^{(t)}$
- Stop when $|r^{(t+1)} - r^{(t)}|_1 < \epsilon$
 - $|x|_1 = \sum_{i=1}^N |x_i|$ is the L_1 norm
 - Can use any other vector norm e.g., Euclidean

10.Dead ends: A page has not out-links \rightarrow cause all pages to 0

11.Solution: Teleport

\rightarrow Allow a surfer to jump to some random page from dead ends

▪ **Power Iteration:**

Set $r_j = 1/N$

1: $r_j' = \sum_{i \rightarrow j} \frac{r_i}{d_i}$

2: $r = r'$

Goto 1

▪ **Example:**

Iteration 0, 1, 2, ...

Iteration 0, 1, 2, ...

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{bmatrix} 1/3 & 1/3 \\ 1/3 & 1/6 \\ 1/3 & 1/6 \end{bmatrix} \begin{bmatrix} 1/3 \\ 1/6 \end{bmatrix} = \begin{bmatrix} 1/4 & 5/24 \\ 1/6 & 1/8 \\ 1/6 & 1/8 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

▪ **Measures generic popularity of a page**

- Ignore or miss topic-specific authorities
- Solution:** Topic-specific PageRank

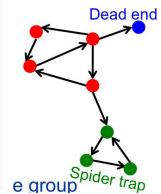
▪ **Uses a single measure of importance**

- Other models of importance
- Solution:** Hubs-and-Authorities

▪ **Susceptible to link spam**

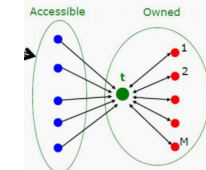
- Artificial link topologies created in order to boost page rank
- Solution:** TrustRank

Limitations of PageRank:



12.Problem: Spider Traps: A group of pages with no links out of the group **13.Solution:Teleport**

15.Topic-Specific PageRank: Teleport set is restricted to a topic specific set: a.Decide on topics b.Pick a teleport set for that topics c. determine the topic (most relevant to a query) d. Use the pageRank



vectors **16.TrustRank:(1) Link Farms:** Maximize the pagerank of page t **(2)TrustRank:** Topic-specific PageRank with a teleport set of trusted pages / Approximation isolation

$$y = x + \beta M \left[\frac{\beta y}{M} + \frac{1-\beta}{N} \right] + \frac{1-\beta}{N} \quad y = \frac{x}{1-\beta^2} + c \frac{M}{N} \quad \text{where } c = \frac{\beta}{1+\beta}$$

(3) Seed pages: make small and point to adequate trust rank

17.Trust Propagation: Each page has trust value (0-1), spam: page below trust threshold.

$$\beta t_p \wedge o_p \quad \text{for } 0 < \beta < 1$$

Trust is additive: trust of p = sum of inlinked pages **Trust is attenuation:** decreases with the distance **Trust is splitting:** split by outlinks

18.Hubs and Authorities: 1.HITS(Hypertext-induced Topic Selection) 2.Authority:contains useful information 3.Hub:link to authorities

- HITS algorithm in vector notation
 - Set: $a_i = h_i = \frac{1}{\sqrt{n}}$
 - Repeat until convergence
 - $h = A \cdot a$
 - $a = A^T \cdot h$
 - Normalize a and h
 - Now $a = A^T \cdot (A \cdot a)$ $h = A \cdot (A^T \cdot h)$
- Convergence criterion:
 - $\sum_i (h_i^{(t)} - h_i^{(t-1)})^2 < \epsilon$
 - $\sum_i (a_i^{(t)} - a_i^{(t-1)})^2 < \epsilon$
- a is updated (in 2 steps):
 - $a = A^T (A a) = (A^T A) a$
 - h is updated (in 2 steps):
 - $h = A (A^T h) = (A A^T) h$

4.PageRank vs HITS: 1.both use link structure to rank page 2.PageRank computed for all web pages and stored prior to query / HITS operates on a small subgraph from web

graph 3.PageRank model: the value of the link depends on the links into u / HITS model: depends on the value of the other links out of u (vulnerable to spam)

Column \rightarrow out-
row \rightarrow in

	y	a	m
y	1/2	1/2	0
a	1/2	0	1
m	0	1/2	0