

UNIT 5

DATA MINING

- It is also known as knowledge discovery (KDD) from data.
- Extracting hidden info or hidden patterns.
- Data mining is also known as KDD or knowledge discovery in database.
- Data processing using sophisticated data search capabilities and statistical algorithms to discover patterns and correlations in large pre-existing databases. A way to discover new meaning in data.
- Supervised and Unsupervised algorithms are used to identify the hidden patterns in data.
- Datasets can be generalized into 3 types
 - i) Data with labels → supervised
 - ii) Data without labels → unsupervised
 - iii) Data with only a small portion of labels.

Supervised algorithm example:

→ Classification

Characteristics of examples:

• Clustering

- * Supervised approaches depend on some a-priori knowledge of the data (classes)
- * Unsupervised algorithms are used to characterize data without any prior information about what kinds of patterns will be discovered by the algorithm.
- * Classification is a common supervised approach and is appropriate when the dataset has labels of a store portion of the data without labels
- * Clustering is a common unsupervised mining technique that is used to find belongingness of datasets without labels
- * for classification - The algorithm learns from the training data and builds a model that will automatically categorize new data elements into one of the existing classes provided with the training data.

- for clustering - Clustering algo determines which elements in the dataset are similar to each other based on the similarity of the data elements.
It makes use of euclidean distance.
- Clustering techniques use keywords that are represented as a vector (To represent a cluster) and the cosine similarity measure is used to distinguish how similar one vector (data elements) is to another

Motivations for Data mining in social media

- The data available via social media can give us insights into social networks or societies that were not previously possible in both scale and extent
- Data mining techniques can help effectively deal with the 3 main challenges of social media data
 - ① Social media datasets are large. Eg: FB
 - ② SM datasets can be noisy & Missing data. dynamic
 - ③ Data from online SM is streaming.

mined but

- Analyzing datamining techniques to large social media datasets has the potential to continue to improve search results for everyday search engines (e.g. google search engine). realize esp specialized target marketing for businesses, help psychologists to study behaviors.
- Provide new insights into social structure for sociologists, personalized web services for consumers and even help detect & prevent spam for all of us.

Data mining methods for social media

Applications that apply data mining techniques developed by industry and academia are already being used commercially

Eg: Somepoint, a social media analysis company, provides services to mining and monitor social media to provide live information about how

Relationships are association b/w individuals, are represented as links in the graph.

A graph representation enables the application of classical representation of mathematical graph theory, traditional social network analysis method, and work on mining graph data.

Representing online social media data as graphs enables leveraging work done with graph related to influence propagation for community detection, and link prediction.

Datamining - A process

Different datasets and data questions require diff type of tools. If it is known how data can be organised
• Verification tool will be appropriate
If you understand what the data is about & cannot associate

friends and followers in clusters, a clustering
clustering tool may be used.

Before determining the raw data to be prepared,
one

With Data represented as a graph the work
begins with the selected number of nodes,
known as seeds. Graphs are traversed
beginning with a set of seeds and
as the link structure from the seed nodes
is exploited, data is collected and
the structure itself is also analyzed. Using
the link structure to extract from the
seed set & gather new information is
known as Crawling the network.

As the crawler discovers new information,
it stores the new information as a
repository for further analysis of cluster.

+ 2 specific types of SM data to

i) Social Networking Sites \rightarrow FB

ii) Blogosphere.

A social networking site like FB / LinkedIn, consist of connected users forming user profile.

Fig: Hypothetical graph structure diagram for a typical social networking site

Social networking site provides excellent source of data for studying collaboration relationships, group structure and who talks to whom.

Thus a graph structure should be utilized

The driving factors for mining social networking sites is the "unique opportunity to understand the impact of a person's position in the network on everything from their tastes to mood to their health".

The most common data mining applications related to social networking site include

- 1) Group Detection
- 2) Finding and identifying the group

Finding individuals that associate more with each other than other users.

A clustering approach such as modularity maximization can be used to detect the subgroups

(d) Group Profiling

"What is this group about?" Advanced data mining techniques such as Topic Taxonomy, a tree structure can be utilized.

(e) Recommendation Systems

A recommendation system analyzes social networking data & recommends new friends or new groups to a user. Eg: Facebook.

- friend suggestion. The ability to recommend group membership to an individual is advantageous for the group that would like to have additional members & can be helpful to an individual who is looking to find other individuals with similar interests or goals.

Recommendations are based on user profile data & a user's association with other users can be used to provide suggestions to users. There are 3 steps for this

- (1) Identify features of the profile.
- (2) Group members are clustered.
- (3) A decision tree is created.

Blogosphere

Need not be blogs

Web logs of blogs or user published journals available on the web.

Blogs are typically open to the public & provide a mechanism for readers to comment on specific post.

The set of all blogs and blog posts is referred to as ~~is~~ blogosphere.

Clustering, matrix factorization and ranking are the most commonly used techniques for data mining in blogosphere.

4 measures / factors to identify influential nodes

(1) Recognition

(2) Activity Generation

(3) Novelty

(4) Eloquence

(5) Topic detection and change

(6) Sentiment analysis

The most common data mining application related to the blogosphere include -

(i) Blog classification

- Organize blogs by topic

Using SVM (Support Vector Machine) automatically blogs can be classified using descriptive words (tags).

- Tags are assigned to the blog post using SVM classifier.

(ii) Identifying influential nodes

A blogger is influential if he/she has the capacity to effect the behaviour of fellow bloggers.

The benefit of being able to identify influential bloggers, blogs, blogpost is that marketing efforts for goods & services could be focused on points of influence that are most likely ^{to gain} ~~to be~~ support for a topic, product or essay.

and started peruse and discussed with a social media. Researches in other platforms have applied text mining algorithm and propagation model to try to develop approaches for better understanding how information moves through the global blogosphere.

Data mining technique can be applied to social media to understand database and to make use of data for research and business purpose. Representative area include community and group detection, information diffusion, influence propagation, topic detection and monitoring, individual behaviour analysis, group behaviour analysis and marketing research for business.

Data Representation

It is common to use a graph representation to study social media datasets. A graph consists of a set containing vertex, node and edges (links) individual or specifically represented as nodes of a graph.