Review of Data Mining Algorithms

Supervised learning

• Classification/prediction

Unsupervised learning

- Clustering
- Association rule mining

Semi-supervised learning

• Active learning

Recommender systems

- Collaborative filtering
- Matrix completion

Graphical models

Terminologies

- Dataset
 Training set
 Testing set
 Validating set
- Data representation □Feature vector
- TF/IDF $idf(this, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$

Document 1		
	Term	Term Count
	this	1
	is	1
	а	2
	sample	1

Term	Term Count
this	1
is	1
another	2
example	3

Document 2

$$\begin{split} \mathrm{tf}(\mathsf{example}, d_2) &= 3\\ \mathrm{idf}(\mathsf{example}, D) &= \log \frac{2}{1} \approx 0.3010\\ \mathrm{tfidf}(\mathsf{example}, d_2) &= \mathrm{tf}(\mathsf{example}, d_2) \times \mathrm{idf}(\mathsf{example}, D) = 3 \times 0.3010 \approx 0.9030\\ {}_{\mathrm{BUDT\,758}} & 2 \end{split}$$

Supervised Learning

- Regression
 Linear regression
 Logistic regression
- Naïve Bayes
 Strong independence assumption
- K-nearest neighboring (KNN)
- Decision Tree
 - **C**4.5

Can handle both numerical and categorical featuresMissing values

Support Vector Machine

• Find a hyper-plane to maximize the functional margin.









5

Unsupervised Learning

- Clustering
 K-means
 Spectral clustering
 Hierarchical clustering
 Density-based clustering (DBSCAN)
 Distance metric
 - EuclideanManhanttanCosine



K-Means



 k initial "means" (in this case k=3) are randomly generated within the data domain (shown in color).



2) k clusters are created by associating every observation with the nearest mean. The partitions here represent the Voronoi diagram generated by the means.



3) The centroid of each of the k clusters becomes the new mean.



4) Steps 2 and 3 are repeated until convergence has been reached.

- Graph-based community detection
 Modularity maximization-based
 - Networks with high modularity have dense connections between the nodes within modules but sparse connections between nodes in different modules.



Semi-supervised learning

• Active learning



Representation learning



Recommender Systems

- User-based collaborative filtering
- Item-based collaborative filtering



Graphical Models: Topic Modeling

Topics

Documents

Topic proportions and assignments

