KDD          DM

DM ⎰ -selecting data set/subset,
    ⎱ - applying algo.,
      - using data for predictive modeling analysis

* DW — data warehousing = implementation of enterprise data in unified structure (cube)
* OLAP — ?.why — faster query response

* MS BI perspective (many data sources, ...)
  * BIDS
* client tools — excel, SP, SSRS, SSIS

DM processing

Training data

$\downarrow$

DM Engine

⊛   * one way

Mining Model } Predictive models

Ⓐ  Mining model    data to be predicted

$\downarrow\downarrow$

data   with predictions     * other way

---

Steps for building DM model
① Model definition (define columns for cases: visually (BIDS), using ~~DMX~~ DMX, or from PMML)

② Model Training ( feed lots of data from a real DB, or from a system log)

③ Model Testing ( testing data must be different from training )

④ Model Use ( exploration and prediction)
  - use the model on new data to predict outcomes

⑤ Model update (monthly, weekly, nightly, ... and re-test)

## Mining structure

~~data~~ – describes data to be mined
 – columns from a data source and their:
  – data type
  – content type

~~contains mining model~~ – contains mining models (often we build several different models in one structure)
 – holds training data, known as cases (if required)
 – holds ~~training~~ testing data, known as Holdout (in SQL 2008)

---

## DM Model

 – container of patterns discovered by DM Algorithm
  amongst the training cases (data)
  – a table containing patterns
   – expressed by visualizers
 – specifies usage of columns already defined in the mining structure

---

### Cases: (the things we study)

 – case – set of columns (attributes) you want to analyse
   ex: age, gender, region, annual spending
 – case key – unique ID of a case
 – a column has:
   * data type
   * content type
   * and optionally:
       * distribution
       * discretization
       * Related columns
       * flags (e.g. NOT NULL)

column distribution ③

- if you know the distribution of your data (you should),
   indicate it:
   - Normal (typical Gaussian bell- ~~~ curve)
   - Log Normal (most values at the "beginning" of the scale
   - Uniform (flat line - equally likely or perfectly random)
- other distributions can exist, but you can not indicate them -
   algorithm will work fine

---

Create Model $M_1$ &lt;name&gt;
            mining
   { Attributes


   } using Algo1

---

data
training

INSERT INTO [Mining model |mining] &lt;name&gt;
                        |structure
      @ [ ( attributes) ]
   &lt; source data &gt;
          → data query / DMX Query / MDX Query
          Stored Procedure call / Rowset
                                 parameter

---

predict

      Select [TOP &lt;count&gt;]
      &lt;expression-list&gt; from &lt;model&gt;
   & [ [Natural] PREDICTION JOIN
          &lt;source data&gt; AS &lt;alias&gt;
          [ON &lt;column- mapping&gt; ]
          [When &lt;filter expression&gt;]
          [order by &lt;expression&gt;]

      ]

DM Algo.

Decision Tree — finds the odds of an outcome based on values in a training set

Association Rules — identifies relationships between cases

Clustering — classifies cases into distinctive groups based on any attribute sets

Naive Bayes — clearly shows the differences in a particular variable for various data elements

sequence clustering — groups or clusters date based on a sequence of previous events

Time series — Analyzes and forecasts time-based data combining the power of ARTXP (developed by my Research) for short-term predictions with ARIMA (in sequence) for long-term accuracy.

Neural Nets — seeks to uncover non-trivial intuitive relationships in data

Linear regression — determines the relationship between columns in order to predict an outcome

Logistic regression — determines the relationships between columns in order to evaluate the probability that a column will contain a specific state.

DM Lift chart for comparing algorithms

logistic regression for risk association rule using time correlation

Scatter plot

## DM

- pattern discovery
- Intelligent grouping
- Predictions (probabilities) (values, series Events, Time-based events)

## Algorithm

- clustering = grouping (~~segmenting~~ divide data ~~into~~ into clusters segments)
- Classification = predicting a specific value
- association = correlation (market basket)
- regression = forecasting a continuous number
- sequences = process and route (clickstream)
- deviation = outliers (exception, fraud)

---

## MS DM

Mining structures (containers to keep DM models)

---

Naive Bayes — simple, starting point
- Basic groupings
- Discrete (unique) content only

Decision Trees - most common, groups
- powerful viewer
- discrete or continuous attributes

Time series — forecasting
- time series viewer

① ⟩

---

Clustering - grouping
sequence clustering - grouping + sequences
~~state tran~~ — state transitions

Association rules - market basket
- itemsets

② ⟩

Neural Network — solves the difficult problems
- finds invisible patterns
- takes significant processing overhead
Linear/logistical regression (Determines best straight line through series of points)

Regression Analysis:=

$$\boxed{\text{Dataset}} \longrightarrow \boxed{\text{RA}} \longrightarrow \boxed{\text{predictive model}}$$

(output) $X =$ predictor variable (output)

$Y =$ response variable (input)

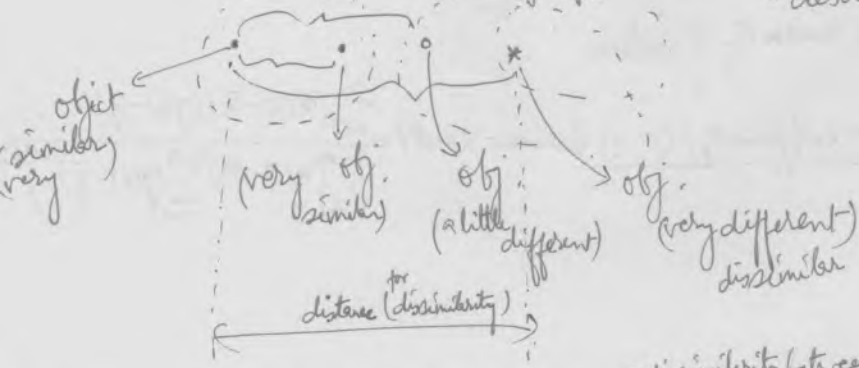① linear function, $Y = aX + b$ $\longrightarrow$ linear (a relation between $X$ and $Y$)

~~(② Score function = $\sum \left( \quad + (\text{...})\right)^2$~~

② Score function = $\sum_{i=1}^{n} \left( y(i) - \hat{y}(i) \right)^2 \rightarrow$ { predicting $n$ "target" values $y(i)$, $1 \leq i \leq n$, predictions for each (squared error)

other score functions { e.g., least squares, classification accuracy, likelyhood, misclassification rate

---

Distance measures:- ( similarity measures between objects } = proximity
~~dissimilarity~~ dissimilarity measures " " }

~~obtain~~
- obtained directly from the objects
- obtained indirectly from vector of measurements/characteristics describing each object.

object (similar) (very)    (very obj. similar)    obj. (a little different)    obj. (very different) dissimilar

distance (for dissimilarity)

① $s(i,j)$, similarity between $i$ and $j$ objects $= 1 - d(i,j)$ $\rightarrow$ dissimilarity between objects $i$ and $j$

② $d(i,j) = \sqrt{2(1-s(i,j))}$

metric (dissimilarity measure) with conditions ① $d(i,j) \geq 0$ and $d(i,j) = 0$ if and only if $i = j$
② $d(i,j) = d(j,i)$ for all $i,j$
③ $d(i,j) \leq d(i,k) + d(k,j)$ for all $i,j,k$ (triangle inequality)

Euclidian distance, $d_E(i,j)$ between $i$ and $j = \left( \sum_{k=1}^{p} (x_k(i) - x_k(j))^2 \right)^{1/2}$    ①!?

   — for $n$ data object with $p$-real valued measurements on each object.

   — for vector of observations for the $i$th object by $X(i) = (x_1(i), x_2(i), \ldots, x_p(i))$,

$$1 \le i \le n$$

❋ there is assumption of some degree of commensurability between different variables.

$\underline{\text{standard} \cancel{\text{dod}} \text{ deviation}}$ (for the $k$th variable, $X_k$) $= \hat{\sigma}_k = \left( \frac{1}{n} \sum_{i=1}^{n} (x_k(i) - \mu_k)^2 \right)^{1/2}$

$\mu_k$ = mean for variable $X_k$

sample mean $\bar{x}_k = \frac{1}{n} \sum_{i=1}^{n} x_k(i)$

$x'_k = \left( x_k / \hat{\sigma}_k \right)$ removes effect of scale as captured by $\hat{\sigma}_k$

weighted Euclidian distance measure, $d_{WE}(i,j)$

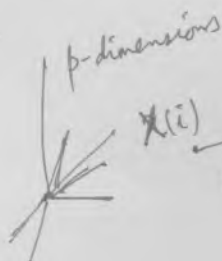$$= \left( \sum_{k=1}^{p} w_k (x_k(i) - x_k(j))^2 \right)^{1/2}$$

$\underline{\text{sample covariance}}$ between $X$ and $Y$

$$Cov(X,Y) = \frac{1}{n} \sum_{i=1}^{n} (x(i) - \bar{x})(y(i) - \bar{y})$$

$\bar{x}$ = sample mean of $X$ values

$\bar{y} =$ "   "   " $Y$ "

$\underline{\text{sample correlation coefficient}}, \rho(X,Y)$ between $X$ and $Y = \dfrac{\sum_{i=1}^{n}(x(i)-\bar{x})(y(i)-\bar{y})}{\left( \sum_{i=1}^{n}(x(i)-\bar{x})^2 \sum_{i=1}^{n}(y(i)-\bar{y})^2 \right)^{1/2}}$

p-dimensions

$X(i)$ —————— $X(j)$

Mahalanobis distance, $d_{MH}(i,j) = \left(X(i) - X(j)\right)^T \Sigma^{-1} \left(X(i) - X(j)\right)$

$T$ = transpose matrix

$\Sigma$ = $p \times p$ sample covariance matrix

$\Sigma^{-1}$ = standardizes data relative to $\Sigma$

✓ Minkowski space
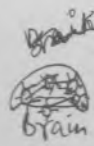✓ $L_\lambda$ metric
✓ Manhattan/city-block metric
✗ Jaccard coefficient
✗ dice coefficient

date = information

small date    large date    pieces of information
a piece of information

example: man is selfish by nature.    → philosopher    brain    data mining
                                                         brain    information

knowledge discovery process    interacted society
contain data mining

implementation logic    algorithm    neurons interact
of    in terms of pseudocode    Biological data mining strike
    or programming code    gathers all the pieces of
                                information and mine
human brain    to discover knowledge them

for large date    not reduced small but efficient miner    Biological limitation
    larger computer RAM    is slow    lower RAM / memory space

cluster
stars
bunches

Data mining?

- finding interesting & structure in data
  - structure: (refers to statistical patterns, predictive models, hidden relationships)

  - interesting: ?

ex: - predictive modeling (classification, regression)
    - segmentation (data clustering)
    - affinity (association) (summarization) - relations between fields, association, visualization

beyond data analysis

  * Scaling analysis to large databases
    - how to deal with data without having to move it out?
    - are the then abstract primitive accesses to data, in database systems, that can provide mining algorithms with the information to drive the search for patterns?

  * automated search
    - enumerate and search. create numerous hypothesis
    - fast search
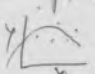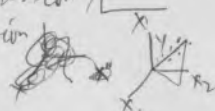    - useful data reductions

(*) more emphasis on understandable models

    ? high dimensions ?

DM

Defining the goals (use models)

① Segmentation — to segment customers by profitability and market potential.

② Profile analysis — analysis of (customer's/prospect's) profile like average age, gender, length of age of customer residence is relationship

③ Response — receiving response of a customer by offering a product/service.

④ Risk — financial risk (grant loan), credit card risk, risk of fraud

⑤ Activation — (predict response of customers, and predict activation given response of customer) by

⑥ Cross-sell and up-sell — (predict the probability or value of a current customer buying a different product/service from same company)
(predict the probability or value of customer buying more of the same product/service)

⑦ Attrition/churn — act of customers switching companies to take advantage of better deal.

⑧ Net Present value — NPV predict overall profitability of a product for a predetermined length of time.

⑨ Lifetime value — LTV predicts overall profit profitability of customer/business for (customer lifetime value) a predetermined length of time.

Modeling methodology (for predictive/descriptive models)

① Linear regression (statistical techniques) — quantifies relationship between two continuous variables.

  ① dependent    ② independent
       $Y$            $X$
                  (predictive variable)

* finding a line, through the data, that minimizes the squared error from each point

* R-square (measuring strength of relationship) — measures the amount % of overall variation in data that is explained by the model %

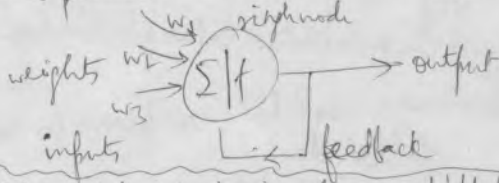∨ Nonlinear (curvilinear) regression
∨ multiple linear regression

## Logistic regression —

* similar to linear regression
* dependent variable is non-continuous (discrete/categorical) logical
* based on statistical distribution

dependent var. → $Y$ purpose

income $X$

## Neural network —

* the process is one of pattern recognition and error minimization.
* made up of nodes that are arranged in layers
* data is split into training and testing data sets. A third group is held out for final validation. Then weights or inputs are assigned to each of the nodes in the first layer. By iterations, the inputs are processed through the system and compared to the actual value. The error is measured and fedback to adjust the weights. The process ends when a predetermined minimum error level is reached.

weights  $w_1$  $w_2$ →  peph node
$w_3$  →  Σ|f  → output
inputs  ← feedback

## Classification tree (decision tree) — to sequentially partition the data to maximize the differences in the dependent variable.)

* classify data into distinct groups/branches that create the strongest separation in the values of the dependent variable.

Associations and item-sets:
rule; if X then Y
$$X \rightarrow Y$$

exception: if others...
if _ and _ then _ except if _
then _
if _ then _ else _

for any rule if $X \rightarrow Y \Rightarrow Y \rightarrow X$ then X and Y are called
on interesting item-set

coverage (predict correctly) = support for a rule $R = \dfrac{\text{no. of occurrences of R}}{\text{total no. of all occurrences of all rules}}$

accuracy = confidence of a rule $X \rightarrow Y = \dfrac{\text{no. of occurrences of Y given X} (X \rightarrow Y)}{\text{no. of all other occurrences given X}}$

for

Apriori Algo: - based on combination
- min. support & min. confidence

Mining association rules using Apriori —
- use Apriori to generate frequent itemsets of different size
at each iteration.
- divide eg each frequent itemset (I) into LHS and RHS
$$LHS \rightarrow RHS$$
- confidence of such rule = support of $\dfrac{\text{itemset}}{(I)}$ / support (LHS)
- discard all rules whose confidence is less than
min confidence.

Classification Techniques : ( predefined classes)
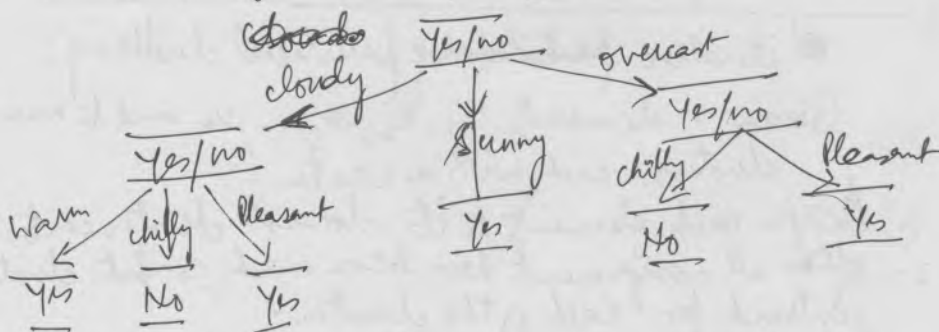decision tree identification :            based on
    (Hunt's Algo. method for .. )              supervised learning)

| weather  | Temp | Play? |
|----------|------|-------|
| Sunny    | 30   | Yes   |
| overcast | 15   | No    |
| Sunny    | 16   | Yes   |
| —        | —    | —     |
| —        | —    | —     |

(set these three into classes)  Warm / chilly / Pleasant

So, we got,     whenever
               Sunny  ⟶  Yes  X
               cloudy  ⟶  Yes/No  } make
               overcast ⟶ Yes/No  }  true

So, by considering three classes,

cloudy        Yes/no        overcast
cloudy                       Yes/no

Yes/no        Sunny        chilly        Pleasant
warm  chilly  Pleasant       Yes          No        Yes
Yes    No     Yes

# Clustering techniques (unsupervised learning)

→ clustering partitions data set into clusters or equivalence classes (but <u>not</u> the <u>predefined classes</u>)

→ similarity among members of a class more than similarity among members across classes.

→ ~~similar~~ similarity measures: Euclidian distance or other application specific measures.

---

⊘ Nearest neighbour clustering Algo:

Given $n$ elements $x_1, x_2, x_3, x_4, x_5, \ldots, x_n$ and threshold $t$.

① $j \leftarrow 1, k \leftarrow 1$, clusters = { }

② Repeat (~~similarity~~)
    ① find nearest neighbour of $x_j$
    ② let nearest neighbour ~~be~~ in cluster $m$
    ③ if distance to nearest neighbour $> t$, then create a new cluster and $k \leftarrow k+1$; else assign $x_j$ to cluster $m$
    ④ $j \leftarrow j+1$

③ until $j > n$

---

⊘ iterative ~~partitioning~~ partitional clustering:

Given $n$ elements $x_1, x_2, x_3 \ldots x_n$ and $k$ number of clusters, each with a center.

1. Assign each element to its closest cluster center
2. After all assignments have been made, compute cluster centroids for each of the cluster.
3. Repeat above two steps with new centroids untill algo. converges.

# Mining streaming data :

example : stock market quotes

**Running mean :**

let $n$ = no. of items read so far

$avg$ = running average calculated so far

on reading the next number $num$ :

$$avg \leftarrow (n * avg + num)/(n+1)$$
$$n \leftarrow n+1$$

---

**Running variance :**

$$var = \sum (num + avg)^2$$
$$= \sum num^2 + 2 * \sum num * avg + \sum avg^2$$

let $A = \sum num^2$ of all numbers read so far

$B = 2 * \sum num * avg$ of all numbers read so far

$C = \sum avg^2$ of all numbers read so far

$avg$ = average of numbers read so far

$n$ = no. of numbers reads so far

$$A \leftarrow A + num^2$$
$$B \leftarrow B + 2 * avg * num$$
$$C \leftarrow C + avg^2$$
$$var = A + B + C$$

Mining streaming data:

$\gamma$ - consistency:

→ Let streaming data be in the form of "frames" ~~events~~ (events)
 where each frame comprises of one or more data
 elements.
 frame (event)

→ support for data element $k$ within a frame is defined
 as (# occurrences of $k$) / (# elements in frame)

→ $\gamma$ - consistency for data element $k$ is the "sustained"
 support for $k$ over all frames read so far, with
 a "leakage" of $(1-\gamma)$



$\gamma *$ support($k$)

leaking rate $(1-\gamma)$

$$level_k(k) = (1-\gamma) * level_{k-1}(k) + \gamma * support(k)$$

## mining sequence data (ordered data)

A sequence is a list of itemsets of finite length.

{pen, pencil, ink} {pencil, ink} {ink, eraser} {ruler, pencil}

(itemsets)

* ★ order of items within an itemset does not matter but
order of itemsets matter

* A subsequence is a sequence with some itemsets deleted.

Let a sequence $S' = \{a_1\} \{a_2\} \{a_3\} \ldots \{a_n\}$ is said to be
contained within another sequence $S$, if

$S$ contains a subsequence $\{b_1\} \{b_2\} \{b_3\} \ldots \{b_m\}$
such that $a_1 \subseteq b_1$, $a_2 \subseteq b_2$, $a_3 \subseteq b_{3n} \ldots$
$$a_n \subseteq b_m$$

Apriori Algo. for sequence data :—
① $L_1 \leftarrow$ set of all interesting 1-sequences
② $k \leftarrow 1$
③ while $L_k$ is not empty do
    ① generate all candidate $k+1$ sequences
    ② $L_{k+1} \leftarrow$ set of all interesting $k+1$ sequences
④ done

mining sequence date ( ordered combination = permutation + self concatenation)

example:

a b c d c
b d a e
a c b d
b e
c a b d a
a a a a
b a c a
c b d b
a b b a b
a b d c

min. support = 0.5

interesting 1-sequences:

a
b
d
e

Candidate 2-sequences:

aa, ab, ad, ae
ba, bb, bd, be
da, db, dd, de
ea, eb, ed, ee

min. support = 0.5
interesting 2-sequences:

ab, bd

Candidate 3-Sequences: (by permuting ab, bd with a,b,d,e 1-sequence)

aba, abb, abd, abe,
aab, bab, dab, eab,
bda, bdb, bdd, bde,
bbd, dbd, ebd

interesting 3-sequences = { }

mining sequence data:

Language interface:

( input set of sequences $\Rightarrow$ output state machine) finite state machine)

"shortest run generalization" Algo. by @Srinivasa & Spiliopoulos 2000)

a a bc b



a a c



aabbc

a a b c

ABD
ABCD ABCD

ABCD

BCD

|  | A | B | C | D |
|---|---|---|---|---|
| 1001 | ink | pen | cheese | bag |
| 1002 | milk | pen | juice | cheese |
| 1003 | milk |  | juice |  |
| 1004 | juice | milk | cheese |  |

(E B F) — ink pen cheese bag
(E B F C) — milk pen juice cheese
(E F) — milk juice
(F E C) — juice milk cheese

min. support = 50%

" confidence = 70%

2/3

$[A] = \frac{1}{4} = 25\%$

$[B] = \frac{2}{4} = 50\%$

$[C] = \frac{3}{4} = 75\%$

$[D] = \frac{1}{4} = 25\%$

$[E] = \frac{3}{4} = 75\%$

$[F] = \frac{3}{4} = 75\%$

$B \to C \quad \frac{2}{2}$

B  BC BD BE BF

BC EF

BC DF

BC EF

BC EF

BC EF

BC EF

BC EF

BC EF

BC EF

BC BF

BC  #50%  100%

BE  25%

BF  25%

CE  50%

CF  50%

EF  $\frac{3}{4}$  75%

BC

CE

CF

EF

BCE

BCF

BCEF

CEF  50%

CEE

$C \to EF$

graph mining patterns: frequent substructures → (how to: ① generate frequent substructure candidate.
for characterizing graph sets,
  discriminating different groups of graphs,
  classifying and clustering graphs,
  building graph indices,
  facilitating similarity search.
~~methods for mining frequent subgraphs~~

basic idea:  Let $V(g)$ be vertex set of graph,
            $E(g)$ " edge set " ",
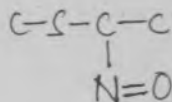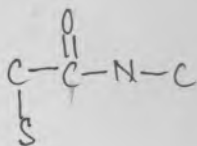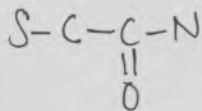   assume $L$ be label function, maps a vertex or edge to label.
   $g' \subseteq g$ ~~subgraph~~
   ($g$, graph is a subgraph of another graph, & $g'$), if there exists
   a subgraph isomorphism from $g$ to $g'$.

$f: V \Rightarrow R^+$

② check frequency of each candidate)

a set of vertices
$$G = (V, E)$$ a set of edges
$$E = \{ \{ (u,v) \mid u, v \in V \} \}$$

a

Given a labeled graph data set, $D = \{ G_1, G_2, \cdots G_n \}$
   the support($g$) or frequency($g$) ∞ is percentage or number
         of graphs in $D$ where $g$ is a subgraph.
A frequent graph is a graph whose support is no less than a
         min. support threshold, min_sup.

Methods for mining frequent subgraphs

① Apriori- based Approach  (search for frequent graphs)
(Apriori Graph)  start with graphs y small size,
         proceed ~~with~~ in bottom-up manner by generating candidates
                        having an extra vertex, edge, path.

S—C—C—N
      ‖
      O

C—C—N—C
  ‖
  S

C—S—C—C
      |
      N=O

① AprioriGraph :

input { D, graph data set
       { min_sup, minimum support threshold

output { $S_k$, frequent substructure set of size k

method : { $S_1 \leftarrow$ frequent single-elements in the data set;
          { Call AprioriGraph (D, min_sup, $S_k$)

procedure : AprioriGraph (D, min_sup, $S_k$)

(1) $S_{k+1} \leftarrow \emptyset$

(2) for each frequent $g_i \in S_k$ do

(3)     for each frequent $g_j \in S_k$ do
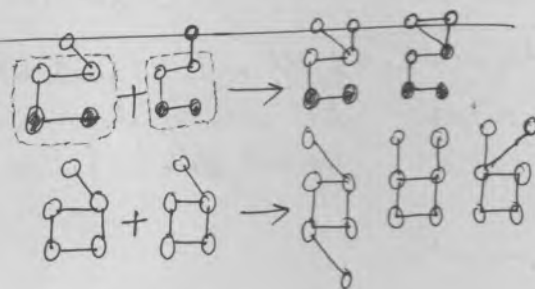
(4)         for each size (k+1) graph g formed by the merge of $g_i$ and $g_j$ do

(5)             if g is frequent in D and $g \notin S_{k+1}$ then
                    insert g into $S_{k+1}$ ;

(6)

(7) if $S_{k+1} \neq \emptyset$ then

(8) AprioriGraph (D, min_sup, $S_{k+1}$)

(9) Return

Algo
{ AGm
{ FSG
{ (edge-disjoint path)
  path-join method

② Pattern-Growth Approach

Input : { g, a frequent graph
         { D, a graph data set
         [ min_sup, minimum support threshold

output: { the frequent graph set, S

Method :

$S \leftarrow \emptyset$

call PatternGrowthGraph $(g, D, min\_sup, S)$

procedure PatternGrowthGraph $(g, D, min\_sup, S)$

① if $g \in S$ then return;
② else insert $g$ into S
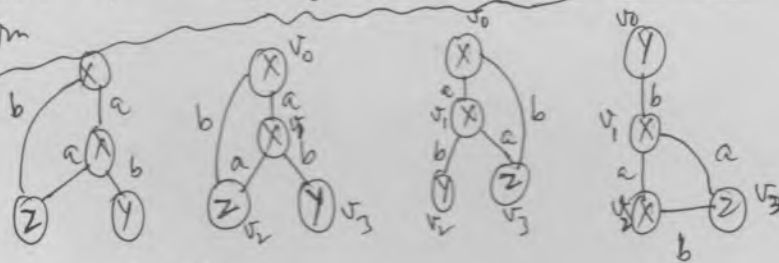③ Scan D once, find all the edges e such that g can be extended
   to $g \diamond_x e$
④ for each frequent $g \diamond_x e$ do {
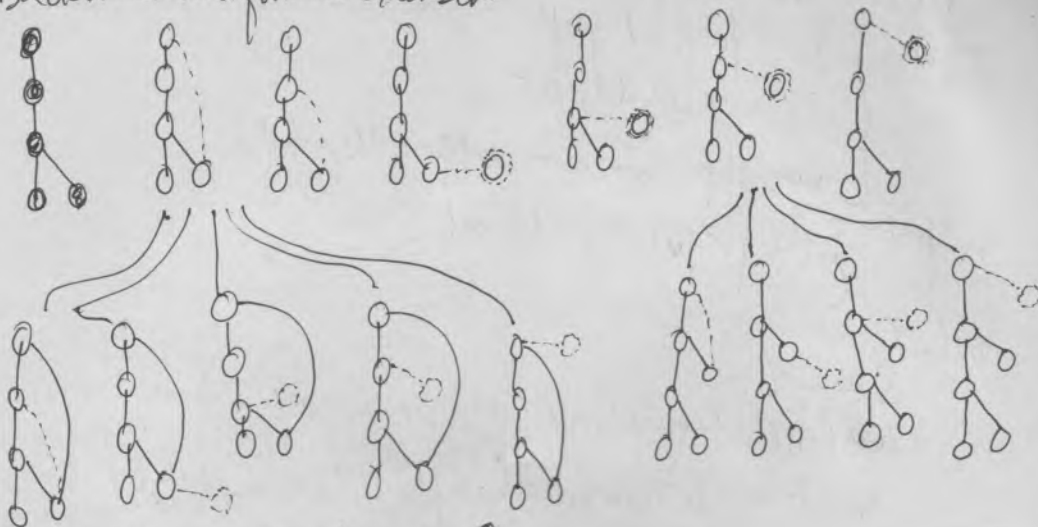⑤ PatternGrowthGraph $(g \diamond_x e, D, min\_sup, S)$ }
⑥ return

gSpan
Algo.



go to the
reduce
generating
duplicate graphs
in
Pattern-GrowthGraph.

DFS subscripting
(Depth First Search)

Backward and forward extension



Right-most extension

# Social Network Analysis (link analysis/ link mining)

— a heterogeneous and multirelational data set represented by graph.

http:// tam. cornell. edu/ strogatz, htm # pub
// www. nd.edu /~ networks/publications. html # talks 0001

x structure always affects function

## characteristics — ( by building graph generated models, then may be used to predict how a network may look in the future)

x if a hypothesis contradicts the generally accepted characteristics, this raises a flag as to the questionable plausibility of the hypothesis. These can help detect abnormalities in existing graphs, which may indicate fraud, spam, or distributed denial of services (DDoS) attacks. Models of graph generation can also be used for simulations when real graphs are excessively large and thus, impossible to collect ( such as very large network friendships).

✓ ① examine the node's degree (no. of edges incident to each node)
② distances between a pair of nodes (shortest path length)
③ network diameter (max. length distance between pairs of nodes)
④ node - to node distances ( the average distance between pairs)
⑤ effective diameter (i.e. the min. distance, d, such that for at least 90% of the reachable node pairs, the path length is at most d)

## Social network phenomena:

① densification power law (number of edges growing superlinearly in the number of nodes) $a \rightarrow$ (1<a<2)

no. of edges; → $e(t) \propto n(t)^a$ ← no. of nodes at time (t)

② Shrinking diameter (the effective diameter tends to decrease as the network grows)

③ Heavy-tailed out-degree "bridging" and in-degree distributions ( the number of out-degrees for a node tends to follow a heavy-tailed distribution by observing the power law $y = \frac{1}{n^a}$ and typically $0<a<2$.

where n = rank of node in the order of decreasing out degrees and typically $0<a<2$.
smaller the value of a, heavier the tail. (Preferential attachment model) "rich get richer"

## Forest Fire model (new nodes attach to the network by "burning" through existing edges in epidemic fashion.)

parameters: forward burning probability, $p$
backward burning ~~probability~~ ratio, $r$

Let a new node '$v$' arrives at time '$t$'. It attaches to $G_t$.

① it chooses an ambassador node, $w$ at random, and forms a link to $w$.

② it ~~selects~~ selects $x$ links incident to $w$, when $x$ is a random number that is binomially distributed with mean $(1-p)^{-1}$. It chooses from out-links and in-links of $w$ but selects in-links with probability $r$ times lower than out-links. Let $w_1, w_2, \cdots \cdots w_x$ denote the nodes at the other end of the selected edges.

③ ~~Our~~ new node, $v$, forms out-links to $w_1, w_2, \cdots; w_x$ and then applies step 2 recursively to each of $w_1, w_2, \cdots; v_x$. Nodes cannot be visited a second time so as to prevent the construction from cycling. The process continues until it dies out.

* (Nodes with heavy-tailed out-degrees may serve as "bridges" that connect formerly ~~so~~ disparate parts of the network, decreasing the ~~so~~ network diameter.

ⓠ $G_i (V, E_i), i = 1, \cdots n$ when $n$ is the number of relations, $V$ is the set of nodes (objects), $E_i$ is the set of edges with respect to the $i$-th relation.

## Link Mining

How can we mine social networks?

Traditional methods of machine learning & data mining, taking as input, a random sample of homogeneous objects from a single relation, may not be appropriate here. So, link mining came in. (it is a confluence of research in social networks, link analysis, hypertext and web mining, graph mining, relational learning, and inductive logic programming. It embodies descriptive and predictive modeling.

* By considering links (the relationships between objects), more information is made available to the mining process. This brings about several new tasks. Here, we list these tasks with examples!
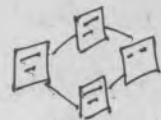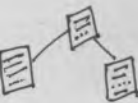
① Link-based object classification (traditionally, objects are classified based on the attributes that describe them.
   Link-based classification predicts the category of an object based not only on its attributes, but also on its links, and on the attributes of linked objects.
   example 1: web page classification (predicts the category of a web page based on word occurrence and hyperlink words/anchor text), both of which serve as attributes.
   ex. 2: In the bibliography domain, objects include papers, authors, institutions, journals and conferences. A classification task is to predict the topic of a paper based on word occurrence, citations, and cocitations (other papers that are cited within the paper), where the citations act as links.
   ex. 3: in epidemiology, predicting the disease type of a patient based on characteristics (e.g. symptoms) of the patient, and on characteristics of other people with whom the patient has been in contact.

② Object type prediction : (predicts the type of an object, based on its attributes and its links, and on the attributes of objects linked to it.
  ex: In the bibliographical domain, we may want to predict the venue type of a publication as either conference, journal, or workshop.
  ex: In the communication domain, a similar task is to predict whether a communication contact is by e-mail, phone call, or mail.

③ Link type prediction : (predicts the type or purpose of a link, based on properties of the objects involved.
  ex: Given epidemiological data, for instance, we may try to predict whether two people who know each other are family members, coworkers, or acquaintances.
  ex: we may want to predict whether there is an advisor- advisee relationship between two coauthors.

④ Predicting link existence ; ( to predict ~~about a link~~ whether a link exists between two objects or not.)
  ex: predicting whether there will be a link between two web pages, and whether a paper will cite another paper.
  ex: In epidemiology, we can try to predict with whom ~~a~~ a patient came in contact.

⑤ Link cardinality estimation: There are two forms of link cardinality estimation.
  (i) we may predict the number of links to an object (· in-link)
   · similarly, the number of out-links can be used to identify web pages that act as hubs, where a hub is one or a set of web pages that point to many authoritative pages of the same topic.
  (ii) the second ~~form~~ form of link cardinality estimation ·(predicts the number of objects reached along a path ~~of~~ from an object. This is important in estimating the number of objects that will be returned by a query.  ( ex: web page)~~web~~

⑥ Object reconciliation : ( to predict whether two objects are, in fact, the same, based on their attributes and links)
  ~~···~~ This is common ~~a~~ task in info. extraction, duplicate elimination, object consolidation, and citation matching / record linkage / identity uncertainty.
  ex: ~~~~ predicting mirror sites.
  ex:  "     "   for two apparent disease strains

⑦ Group detection : ( predicts whether a set of objects belong to the ~~same~~ same group or cluster, based on their attributes and links.
(clustering)
  ex: identification of web communities

⑧ Subgraph detection: (subgraph identification finds characteristic subgraphs within networks.) This is a form of graph search.
ex: discovery subgraphs corresponding to protein structures.

⑨ Metadata mining: (the metadata provide semi-structured data about unstructured data, ranging from text and web data to multimedia databases)
* useful for data integration tasks in many domains
ex: scheme mapping; scheme discovery; scheme reformulation.

---

* <u>Link prediction</u>: what edges will be added to the network?
approaches to link prediction have been proposed based on several measures for analyzing the "proximity" of nodes in a network.
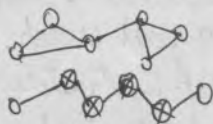
general methodology: All methods assign a connection weight, $score(X, Y)$ to pairs of nodes, $X$ and $Y$, based on the given proximity measure and input graph, $G$. A ranked list in decreasing order of $score(X,Y)$ is produced. This gives the prediction predicted new links in decreasing order of confidence. The predictions can be evaluated based on real observations on experimental data sets.

the simplest approach ranks pairs $\langle X,Y \rangle$, by length of their shortest path in $G$. This embodies the small world notion that all individuals are linked through short chains. (Since the convention is to rank all pairs in order of decreasing score, here, $score(X,Y)$ is defined as the negative of the shortest path length.) Several measures use neighborhood info. The simplest such measure is "common neighbors". The greater the number of neighbors that X and Y have in common, the more likely X and Y are to form a link in the future.
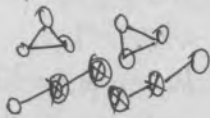
---

Multirelational social network analysis:
(same kind of relation)
* homogeneous and heterogeneous links
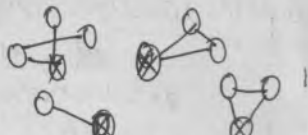different kind of relation
* relation selection and extraction

ⓐ    ⓑ    ⓒ

Let a user requires that the four ⊗ objects belong to the same community and specifies this with a query.

- As for graphs, the relative importance of each of the three relations differs with respect to the user's information need. ⓐ is the most relevant to the user's need, ⓑ comes in second, ⓒ is noisy or negative in regards to the user's information need.

But different aspect of a user can vary the outcome/result.

* In multi-relational social network, community mining should be dependent on the user's query (or information need). A user's query can be very flexible.
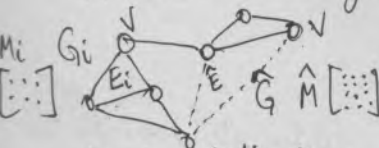
An algorithm for relation extraction and selection was proposed, which models the task as an optimization problem.

for a Given V= a set of objects and ~~a~~ set of relations,
a set of graphs, $G_i (V, E_i)$, ~~a~~ $i = 1, 2, \cdots, n$  (n= no. of relations)
$E_i$ = set of edges (relations) with respect to i-th relation.

The weights on the edges can be naturally defined according to the relation strength of two objects. The algorithm characterizes each relation by a graph with a weight matrix. Let $M_i$ denote the weight matrix associated with $G_i$, $i = 1, \cdots, n$. Each element in the matrix reflects ~~the~~ the relation strength between a pair of objects in the relation.

suppose a hidden relation is represented by a graph $\hat{G}(V, \hat{E})$, and $\hat{M}$ denotes the weight matrix associated with $\hat{G}$. A user specifies her info. need as a query in the form of a set of labeled objects $X = [x_1, \cdots, x_m]$ and $Y = [y_1, \cdots, y_m]$ when   (when $y_j$ is the

$M_i$  $G_i$  $\hat{E}$  $\hat{G}$ $\hat{M}$  $\begin{bmatrix} \cdot \cdot \cdot \end{bmatrix}$

the algo. aims at finding a linear combination of these weight matrices to approximate $\hat{G}$

such labeled objects indicate partial info. of hidden relation → label of $x_j$.

http:// scholar.google.com/
http:// opac.dl.itc.u-tokyo.ac.jp/
http:// www.dl.itc.u-tokyo.ac.jp/gacos/

encyclopedic and dictionary: http://na.jkn21.com/
index to resource: http://resource.lib.u-tokyo.ac.jp/iri/url_search.cgi
                   http://webcat.nii.ac.jp/webcat_eng.html
ebook & materials:    http://webcatplus.nii.ac.jp/
                   http://opac.ndl.go.jp/
         http:// u-tokyo.navi.littel.jp/
         http://ci.nii.ac.jp/en
UT Article:       www.lib.u-tokyo.ac.jp/ext/utas/
web of Science:   http://isiknowledge.com/WOS
Science direct:        www.sciencedirect.com/
Springer link:         www.springerlink.com/
   Wiley Interscience:    www3.interscience.wiley.com/cgi-bin
                                                    /home

engineering village: www.engineeringvillage.org/
   UT dissertation: http://gakui.dl.itc.u-tokyo.ac.jp/
doctoral dissertation:    http://dbr.nii.ac.jp/

e-journal: http:// www.lib.u-tokyo.ac.jp/ext/ejportal/
          http:// ejournal.dl.itc.u-tokyo.ac.jp/

ebook:    http:// www.netlibrary.org/
UT Repository:   http://repository.dl.itc.u-tokyo.ac.jp/
*journal citation report:   http:// isiknowledge.com/JCR
         ↳ impact factor
* Searching Ulrichsweb:    http:// www.ulrichsweb.com/
                   ↳ search journals, bibliographical info. on serials
                                                    published

(DB) SW repository

- isfsg. org
- promisedata. org
- icu-project. org /repository /
- nasa (Nasa metric data program)
  Nasa IV & V facility
- google code (+time line, - - · )
- kdnuggets. com →

DB SW
repository
promisedata
SW rep.

isbsg.org  N

promisedata.org   (85 datasets)

icu-project.org /repository/

nasa    metrics data program
         nasa  IV&V  data facility

google code

kdnuggets.com

<u>Pajek</u> Jung Guess Prefuse . netmap RSNA

Sept 1991, Scientific American, (Communicating computers and Networks), Special issue

    * virtual community

  * (directed graph)

  (newsgroup) usenet → Pajek

~~high~~ degree (~~high number of~~ connections to (a node)

     node (a large

~~indegree~~ (an ~~author~~/node ~~receives~~ (a number of ~~replies~~/ connections) ~~receives~~

~~out-degree~~ (" " " sends " " " " " )

Degree of node (number of connections to ~~the~~ the node)

in-degree ( the number of replies/connection received by a node/author)

out-degree ( the number of replies/connections sent by a node/author)

structural equivalent (same kind of structural ~~patterns~~ as others)

                           connections

* distinguishing attributes (in context with newsgroup)

         — Answer person (* outward ties to local isolates)

                       (* relative absence of triangles)

                       (* Few intense ties )

         — Reply Magnet (* ties from local isolates often inward only )

                    (* sparse, few triangles )

                    (* Few intense ties )