# CS 591.03
# Introduction to Data Mining
# Instructor: Abdullah Mueen

## LECTURE 9: GRAPH MINING

# Graph Similarity

Edit distance/graph isomorphism:
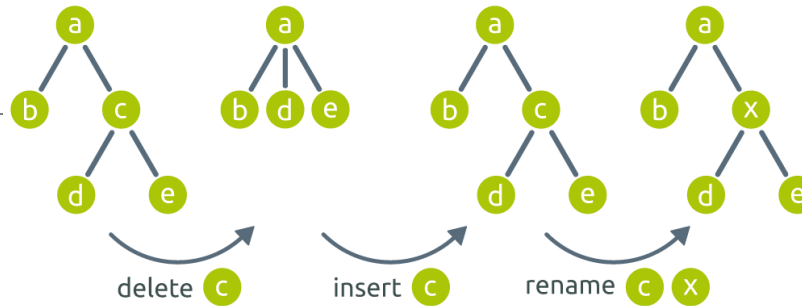◦ Tree Edit Distance

Feature extraction
◦ IN/out degree
◦ Diameter

Iterative methods
◦ SimRank

# Edit Distance

Three Operations

Tai's algorithm runs in $O(m^3n^3)$ time and space for trees with $m$ and $n$ nodes respectively.

# Diameter

Largest Shortest path in the graph.

$$\text{shortestPath}(i, j, 0) = w(i, j)$$

$$\text{shortestPath}(i, j, k+1) = \min(\text{shortestPath}(i, j, k), \text{shortestPath}(i, k+1, k) + \text{shortestPath}(k+1, j, k))$$

```
1 let dist be a |V| × |V| array of minimum distances initialized
to ∞ (infinity)
2 for each vertex v
3    dist[v][v] ← 0
4 for each edge (u,v)
5    dist[u][v] ← w(u,v)  // the weight of the edge (u,v)
6 for k from 1 to |V|
7    for i from 1 to |V|
8       for j from 1 to |V|
9          if dist[i][j] > dist[i][k] + dist[k][j]
10             dist[i][j] ← dist[i][k] + dist[k][j]
11          end if
```

# Simrank

$$s(a, b) = \frac{C}{|I(a)||I(b)|} \sum_{i=1}^{|I(a)|} \sum_{j=1}^{|I(b)|} s(I_i(a), I_j(b))$$

For a node v in a graph, we denote by I(v) and O(v) the set of in-neighbors and out-neighbors of v, respectively.

1. A solution s($*$, $*$) ∈ [0, 1] to the $n^2$ SimRank equations always exists and is unique.
2. Symmetric
3. Reflexive

# Ranking Nodes

## Page Rank

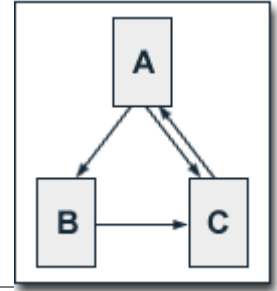$$PR(A) = (1-d) + d (PR(T1)/C(T1) + ... + PR(Tn)/C(Tn))$$

PR(A) is the PageRank of page A,

PR(Ti) is the PageRank of pages Ti which link to page A,

C(Ti) is the number of outbound links on page Ti and

d is a damping factor which can be set between 0 and 1.

# Example



PR(A) = 0.5 + 0.5 PR(C)
PR(B) = 0.5 + 0.5 (PR(A) / 2)
PR(C) = 0.5 + 0.5 (PR(A) / 2 + PR(B))

These equations can easily be solved. We get the following PageRank values for the single pages:

PR(A) = 14/13 = 1.07692308
PR(B) = 10/13 = 0.76923077
PR(C) = 15/13 = 1.15384615

# HITS: Hyperlink-Induced Topic Search

Iterate($G$,$k$)

    $G$: a collection of $n$ linked pages

    $k$: a natural number

    Let $z$ denote the vector $(1, 1, 1, \ldots, 1) \in \mathbf{R}^n$.

    Set $x_0 := z$.

    Set $y_0 := z$.

    For $i = 1, 2, \ldots, k$

        Apply the $\mathcal{I}$ operation to $(x_{i-1}, y_{i-1})$, obtaining new $x$-weights $x_i'$.

        Apply the $\mathcal{O}$ operation to $(x_i', y_{i-1})$, obtaining new $y$-weights $y_i'$.

        Normalize $x_i'$, obtaining $x_i$.

        Normalize $y_i'$, obtaining $y_i$.

    End

    Return $(x_k, y_k)$.

http://www.cs.cornell.edu/home/kleinber/auth.pdf