

Dibujado de grafos mediante Stress Majorization

Oscar Esau Peralta Rosales

12/17/2019

Proyecto final - Métodos Numéricos

El principal objetivo en el campo de Visualización de Datos es ayudar a descubrir, interpretar y comparar relaciones e información existente en un conjunto de datos a través de distintos medios visuales, y uno de estos medios que se analizan en este proyecto son los grafos.

Gran variedad de características presentes en los datos pueden ser modeladas a través de grafos, y gracias a su enorme versatilidad diversos campos de estudio hacen de su uso; geografía, ciencias sociales, movilidad, telefonía, bases de datos, inteligencia artificial, medicina, electrónica, química, etc. Así, a través del proceso de *Stress Majorization* se buscó la una forma de poder dibujar este modelo tal que podamos identificar más fácilmente las relaciones en él conservando la estructura intrínseca del grafo y las similitudes existentes entre los nodos.

Introducción

Un grafo $G = V, E$ se define como un un conjunto de vértices V y aristas E . en dónde cada nodo mapea alguna entidad o característica de los datos y las aristas las relaciones entre estos. Durante este trabajo nos concentramos a trabajar sobre grafos no dirigidos, es decir, las relación presente por una arista es bidireccional.

Dadas las características modeladas de los datos a través de los grafos, podemos intuir que cada grafo es único y por tanto la forma de como dibujarlo de tal manera que se preserve la estructura de sus relaciones también es variada. Diversas técnicas se han realizado con este fin como, *Stochastic Gradient Descent*, *gradient projection*, *Weighted Constraint Relaxation*, *Multidimensional Scaling* y el método implementado aquí *Stress Majorization*.

Stress Majorization es una técnica de optimización usada comúnmente para desplegar información almacenada en una matrix de distancias, la cual contiene la distancia entre cada par de nodos. Esta técnica trata de colocar cada nodo dentro de un *layout* multidimensional tal que la distancias sean preservadas de la mejor manera a través de la función de coste *stress*. La función de coste se define como:

$$stress(X) = \sum_{i < j} w_{ij} (||X_i - X_j|| - d_{ij})^2 \quad (1)$$

Dónde X es un *layout* X de dimensiones $N \times R$ con N como el número de nodos y R la dimensionalidad de este. d_{ij} es la distancia entre el nodo i y el nodo j y w_{ij} es una constante de normalización igual a $d_{ij}^{-\alpha}$.

Así, dado un *layout* inicial X se busca minimizar la función de *stress* para encontrar un nuevo *layout* tal que se tenga un dibujado óptimo.

##Stress Majorization

Dado un *layout* X de dimensiones $N \times R$ con $R = 2$ o $R = 3$ se busca minimizar la función de *stress*

$$stress(X) = \sum_{i < j} w_{ij} (||X_i - X_j|| - d_{ij})^2$$

Expandiendo la función de *stress* se tiene

$$stress(X) = \sum_{i < j} w_{ij} d_{ij}^2 + \sum_{i < j} w_{ij} ||X_i - X_j||^2 - 2 \sum_{i < j} \delta_{ij} ||X_i - X_j|| \quad (2)$$

Con $\delta_{ij} = w_{ij}d_{ij}$ y $w_{ij} = d_{ij}^{-\alpha}$ con $\alpha = 2$ (proporciona mejores resultados).

El segundo término de (2) se puede escribir como

$$\sum_{i < j} w_{ij} \|X_i - X_j\|^2 = (X^T L^w X)^T \quad (3)$$

Donde L^w es la forma cuadrática de la Matriz Laplaciana Ponderada de dimensiones $N \times N$

$$L_{ij}^w = \begin{cases} -w_{ij} & i \neq j \\ \sum_{k \neq i} w_{ik} & i = k \end{cases} \quad (4)$$

El tercer término puede ser limitado, usando la desigualdad de *Cauchy-Schwartz*

$$\|x\| \|y\| \geq x^T y$$

y adicionalmente dada una matriz Z de las mismas dimensiones de X tenemos

$$\|X_i - X_j\| \|Z_i - Z_j\| \geq (X_i - X_j)^T (Z_i - Z_j) \quad (5)$$

$$\sum_{i < j} \delta_{ij} \|X_i - X_j\| \geq \sum_{i < j} \delta_{ij} (\|Z_i - Z_j\|)^{-1} (X_i - X_j)^T (Z_i - Z_j) \quad (6)$$

La ecuación (6) puede ser reescrita de forma matricial como

$$\sum_{i < j} \delta_{ij} \|X_i - X_j\| \geq (X^T L^Z Z)^T \quad (7)$$

Dónde dada una matriz Z de posiciones de las mismas dimensiones de X podemos construir L^Z como

$$L_{ij}^Z = \begin{cases} -\delta_{ij} (\|Z_i - Z_j\|)^{-1} & i \neq j \\ \sum_{k \neq i} L_{ik}^Z & i = k \end{cases} \quad (8)$$

De esta forma la función *stress* puede ser limitada usando las ecuaciones anteriores

$$stress(X) \leq \sum_{i < j} w_{ij} d_{ij} + (X^T L^w X)^T + 2(X^T L^Z Z)^T \quad (9)$$

Diferenciando con respecto a X se obtiene que el mínimo de la función de stress (9) está dado por

$$L^w X = L^Z Z \quad (10)$$

equivalentemente resolviendo para cada eje de X

$$L^w X^{(a)} = L^Z Z^{(a)} \quad a = 1, 2, 3 \dots R \quad (11)$$

La ecuación (11) puede formularse de manera iterativa para un *layout* inicial $X(t)$ y obtener un nuevo *layout* $X(t+1)$ tal que satisfaga $stress(X(t+1)) \leq stress(X(t))$

$$L^w X(t+1)^{(a)} = L^{X(t)} X(t)^{(a)} \quad a = 1, 2, 3 \dots R \quad (12)$$

El proceso iterativo se detiene mediante la siguiente condición

$$\frac{\text{stress}(X(t)) - \text{stress}(X(t+1))}{\text{stress}(X(t))} \leq \epsilon \quad (13)$$

con $\epsilon \sim 10e^{-4}$. Nótese que L^w es constante durante todo el proceso y $L^{X(t)}$ depende del layout actual.

Implementación y Resultados

Datos

Los conjuntos de datos fueron obtenidos desde la página de The SuiteSparse Matrix Collection elegidos por cantidad de nodos y distribución de los nodos.

Table 1: Datasets usados para dibujado

Nombre	Nodos	Tipo
1138_bus	1138	Power Network Problem
cage9	3534	Directed Weighted Graph
CSphd	1882	Directed Graph
dwt_1005	1005	Structural Problem
dwt_2680	2680	Structural Problem
g_26	2000	Undirected Random Graph
gre_216a	216	Directed Weighted Graph
gre_1107	1107	Directed Weighted Graph
qh882	882	Power Network Problem

Generación del layout inicial

El *layout* solo cuenta con dos dimensiones por tanto la matriz X es de dimensión $N \times 2$. La primer columna de la matriz representa el eje x y la segunda columna el eje y. Así el i -ésimo nodo está posicionado en las coordenadas $(X[i][0], X[i][1])$.

Los valores de las coordenadas iniciales para cada nodo son tomados aleatoriamente y posteriormente cada columna es normalizada tener en la menor manera posible errores numéricos.

Cálculo de las distancias entre los nodos

En la literatura sobre diversas implementaciones se observa que el cálculo de la distancia de los nodos lo realizan a través del *Algoritmo de Dijkstra*, el cual nos ayuda a obtener la distancia de un nodo hacia los demás. Sin embargo, necesitamos la distancia de todos los nodos hacia todos los nodos y así poder construir una matriz de distancias D , donde la intersección de la fila i con la columna j representa la distancia que hay entre los dos nodos i y j se procedió a implementar el algoritmo de *Floyd Warshall*. el cuál hace uso de de la aproximación *Bottom-Up* de programación dinámica para realizar el cálculo de todas estas distancias de unas sola vez.