

Обзор по теме курсового проекта

Поиск сообществ в сложных сетях

Введение

Сложные сети – это графы, обладающие определёнными свойствами, какими обладают большинство графов, встречаемых в реальной жизни: социальные графы, нейронные сети, веб-графы и т.п. Их характерные черты – малый диаметр (как правила порядка $O(\log n)$, где n – количество вершин в графе), большой коэффициент кластеризации, иерархическая структура. Ещё одной характерной чертой является степенной закон распределения различных величин, например степеней вершин, что делает данные сети независимыми от масштаба (scale-free)[1].

Сообществами называются такие подмножества вершин, у которых много рёбер внутри подмножества и мало – с остальными вершинами[2, 3]. Существуют более строгие определения сообщества, но они не общеприняты. Для поиска сообществ существует множество алгоритмов, которые будут рассмотрены в основной части.

Помимо рёбер на основе явно заданных связей между вершинами в сложной сети можно строить рёбра на основе метаданных, ассоциированных с вершинами, соединяя вершины со схожими метаданными. Такой граф называется графом данных (proximity graph). Существует несколько методов построения графа данных[6] (relative neighborhood graph, gabriel graph и т.д.), которые будут рассмотрены в основной части обзора.

Обзор существующих решений

Поскольку сложные сети как правило имеют большое количество вершин и рёбер, масштабируемость используемых алгоритмов критически важна.

Поиск сообществ

Будем рассматривать задачу разбиения вершин графа на непересекающиеся подмножества – сообщества, отметив что помимо неё существует также задача разбиения на пересекающиеся сообщества (soft clustering).

Также не будем рассматривать методы, в которых требуется заранее определить число компонент, среди которых максимизация правдоподобия статистической модели, спектральный метод и другие. Вместо этого рассмотрим лишь алгоритмы, в которых количество сообществ автоматически определяется алгоритмом, исходя из данных.

Важнейший класс таких методов основан на введении метрики качества разбиения вершин на сообщества и оптимизации введённого функционала. В качестве такого функционала часто берут модулярность, определяемую следующим образом.

$$Q = \frac{1}{m} \sum_{i,j \in V} (A_{ij} - P_{ij})[c_i = c_j]$$

Здесь m – количество ребёр, V – множество вершин, A – матрица смежности графа, c_i – сообщество, к которому отнесена вершина i . P_{ij} – некоторая оценка A_{ij} для случайного графа, например $\frac{k_i k_j}{2m}$, где k_i – степень вершины i .

Задача оптимизации модулярности NP-полна, поэтому используемые на практике алгоритмы не ищут точное решение, а лишь стараются найти как можно лучшее приближение. Один из таких алгоритмов, CNM[4] – работает за $O(n \log^2 n)$ для сложных сетей. Он создаёт по одному сообществу на каждую вершину, а затем на каждом шаге объединяет два таких сообщества, что при этом увеличение модулярности будет на данном шаге максимальным. Такие итерации продолжаются. Таким образом, он жадно максимизирует функционал модулярности.

Другой подход к поиску сообществ, label propagation[5], основан на запуске случайного процесса, в результате работы получается разбиение множества вершин на сообщества. Вначале каждой вершине присваивается

своё сообщество, а затем на каждой итерации для каждой вершины считается число её соседей из каждого сообщества и выбирается сообщество с максимальным числом соседей. В зависимости от модификации либо сообщество вершины обновляется сразу после подсчёта (асинхронный вариант), либо сначала для всех вершин считаются выбираются новые сообщества, а потом они одновременно обновляются (синхронный вариант), либо сообщество вершины обновляется сразу, но вершины перебираются в таком порядке, что вершины из одного сообщества идут подряд (полу-синхронный вариант). Итерации продолжаются до тех пор, пока количество поменявших сообщество вершин на каждом шаге не меньше некоторого числа, выбираемого в зависимости от данных.

Каждая итерация имеет временную сложность $O(n)$, при этом число итераций, необходимых, чтобы процесс сошёлся, невелико даже для больших n . Таким образом, данный алгоритм является масштабируемым.

Отметим также метод consensus clustering, когда на основе нескольких разбиений на сообщества (например, сгенерированных одним алгоритмом, зависящим от случайности) строится новое разбиение. Для каждого ребра считается, в сколько из разбиений оно соединяет вершины из одного сообщества. Далее из графа выбрасываются рёбра, для которых это число меньше определённой заданной границы, и в оставшемся графе находятся связные компоненты, которые и будут сообществами.

Построение графа данных

Рассмотрим несколько видов графов данных. Все они используют представление метаданных вершины как вектора вещественных чисел заданной размерности. Поскольку во всех из них используется функция расстояния, то меняя метрику в пространстве, можно получать различные графы.

1. ε -шары. Каждая вершина соединяется с вершинами, на расстоянии не более ε от неё. Для ускорения нахождения такого множества вершин можно использовать различные структуры данных, например R-дерево.
2. Граф ближайших соседей. Каждая вершина соединяется с k ближайшими соседями. Для быстрого нахождения ближайших соседей могут использоваться такие структуры данных, как k-d tree или ball tree.
3. Sphere of Influence Graph. Для каждой вершины строится сфера с центром в этой вершине и радиусом, равным её расстоянию до ближайшего соседа. Две точки соединяются ребром, если соответствующие им сферы пересекаются (включая или исключая касание).
4. Минимальное остовное дерево. Рассматривается граф, в котором между каждыми двумя вершинами проведено ребро с весом, равным расстоянию между этими вершинами, и в таком графе находится минимальное остовное дерево.
5. Relative neighbourhood graph. Для каждой пары точек рассмотрим пересечение двух сфер с центрами в этих точках и радиусами, равными расстоянию между ними. Если в этом множестве нет других точек, то данная пара точек соединяется ребром.
6. Граф Габриеля. Для каждой пары точек рассмотрим сферу, построенную на отрезке между ними, как диаметре. Если в этом множестве нет других точек, то данная пара точек соединяется ребром.
7. Триангуляция Делоне. В d -мерном пространстве это такое множество точек, что для каждого симплекса на d вершинах в его описанной гиперсфере нет других точек.

Заметим, что наивный алгоритм вычисления RNG и графа Габриеля занимает $O(n^3)$, что может быть слишком долго, так как, как упомянуто выше, n обычно принимает большие значения. Существует более оптимальный алгоритм построения RNG, использующий $O(n^{2.5})$ времени, однако его недостатком является количество требуемой памяти – $O(n^2)$. Также для частного случая, а именно метрики L_{inf} существует алгоритм с временной сложностью $O(n^2 \log n)$ [8]. Существуют также быстрые алгоритмы для построения графа Габриеля при размерности векторов 2 или 3 [9].

Заключение

Таким образом, для поиска сообществ в сложных сетях наиболее подходят алгоритмы CNM и алгоритм распространения меток (label propagation).

Список литературы

- [1] M. E. J. Newman, "The structure and function of complex networks", 2003
- [2] Santo Fortunato, Darko Hric, "Community detection in networks: A user guide", 2016
- [3] Santo Fortunato, "Community detection in graphs", 2009
- [4] Aaron Clauset, M. E. J. Newman, and Cristopher Moore, "Finding community structure in very large networks", 2004.
- [5] Usha Nandini Raghavan, Reka Albert, Soundar Kumara , "Near linear time algorithm to detect community structures in large-scale networks", 2009
- [6] Matthias Beck, "Computational Discrete Geometry"
- [7] Jyrki Katajainen, Olli Nevalainen, An almost naive algorithm for finding relative neighbourhood graphs in L_p metrics, 1987
- [8] Joseph O'Rourke, Computing the relative neighborhood graph in L_1 and L_{\inf} metrics, 1981.
- [9] Guiseppe Liotta, Low degree algorithms for computing and checking Gabriel graphs