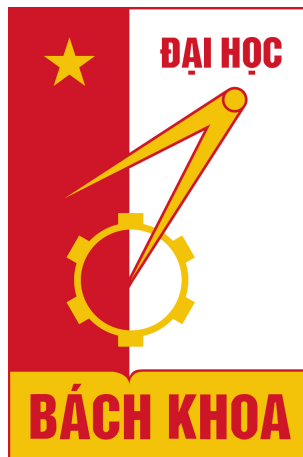


TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN TOÁN ỨNG DỤNG VÀ TIN HỌC

—o0o—



BÁO CÁO ĐỒ ÁN I

**LÝ THUYẾT ĐỒ THỊ VÀ ỨNG DỤNG TRONG
MẠNG LƯỚI GIAO THÔNG**

Học phần: Đồ án I

Giảng viên hướng dẫn: **TS. ĐOÀN DUY TRUNG**

Sinh viên: **Nguyễn Công Hiếu - 20195016**

Lớp: **Toán tin 02 - khóa 64**

HÀ NỘI, 06/2022

Lời nói đầu

Tóm tắt

Trong toán học, lý thuyết đồ thị là một lĩnh vực nghiên cứu về đồ thị - một cấu trúc được sử dụng để mô hình hóa mối quan hệ giữa từng cặp đối tượng. Đồ thị đóng một vai trò quan trọng trong nhiều lĩnh vực nghiên cứu của các bộ môn khoa học. Chẳng hạn như trong khoa học máy tính, các mạng lưới truyền thông, các thiết bị tính toán hiệu năng cao đều được biểu diễn dưới dạng đồ thị. Hay trong vật lý và hóa học, đồ thị được sử dụng để nghiên cứu về mô hình tự nhiên của một phân tử.

Bài báo cáo này được chia thành 3 chương chính:

Chương 1: Cơ sở về lý thuyết đồ thị

Phần này sẽ đi vào những kiến thức căn bản về đồ thị như khái niệm, định nghĩa, tính chất, ... để cung cấp cho người đọc một góc nhìn tổng quan về lý thuyết đồ thị. Một vài định lý, mệnh đề hay hệ quả quan trọng sẽ được đưa ra cùng với cách thức chứng minh để giúp người đọc rèn luyện một tư duy nhạy bén trong các vấn đề về đồ thị.

Chương 2: Trực quan hóa đồ thị với NetworkX API

Với mục tiêu ứng dụng lý thuyết đồ thị vào trong các bài toán thực tiễn, ta cần một phương thức để biểu diễn và hình dung nó trên máy tính. Trong bài báo cáo này, chúng em xin phép được chọn ngôn ngữ python và thư viện NetworkX để giải quyết vấn đề nêu trên.

Chương 3: Ứng dụng trong mạng lưới giao thông của Uber

Sau khi nắm được nền tảng lý thuyết và cách thức xử lý đồ thị trên máy tính, ta sẽ ứng dụng chúng vào phân tích bộ dữ liệu mạng lưới giao thông của Uber. Ở phần này, chúng em sẽ tự đặt ra các vấn đề và tự trả lời chúng dựa trên các kiến thức và thuật toán đã được nghiên cứu ở chương 1 và chương 2.

Mục lục

Lời nói đầu	i
Tóm tắt	ii
1 Cơ sở về lý thuyết đồ thị	1
1.1 Đồ thị và một số khái niệm cơ bản	1
1.1.1 Đồ thị vô hướng	1
1.1.2 Phép đẳng cấu trong đồ thị	5
1.1.3 Phân hoạch trên đồ thị	6
1.1.4 Một số dạng đặc biệt của đồ thị	6
1.1.5 Đồ thị có hướng	6
1.2 Các cách biểu diễn đồ thị	6
1.2.1 Biểu diễn tĩnh	6
1.2.2 Biểu diễn động - danh sách kề	7
1.3 Đồ thị dạng cây	7
1.3.1 Cây và các tính chất cơ bản	7
1.3.2 Cây bao trùm	7
1.4 Đường đi ngắn nhất	7
1.4.1 Thuật toán BFS	7
1.4.2 Thuật toán Dijkstra	7
1.4.3 Thuật toán Floyd-Warshall	7
1.4.4 Thuật toán Bellman-Ford	7
1.5 Mạng lưới và tính liên thông	7
1.5.1 Luồng trong mạng	8
1.5.2 Luồng cực đại	8
2 Trực quan hóa đồ thị với NetworkX API	11
2.1 Giới thiệu về NetworkX API	11
2.1.1 Mô hình dữ liệu	11

2.1.2	Các thống kê cơ bản của một đồ thị	12
2.1.3	Thao tác trên đồ thị	14
2.2	Trực quan hóa đồ thị	15
2.2.1	Hairball	15
2.2.2	Matrix plot	16
2.2.3	Arc plot	17
2.2.4	Circos plot	18
2.2.5	Hive plot	19
3	Ứng dụng trong mạng lưới giao thông của Uber	21
3.1	Xây dựng mô hình đồ thị	22
3.2	Bài toán người du lịch	24
3.3	Khái quát hóa bản đồ giao thông	27
3.4	Lưu lượng giao thông	29
3.5	Bài toán về luồng cực đại	32
3.6	Cắt tỉa trên đồ thị	33
	Kết luận	34
	Tài liệu tham khảo	35

Chương 1

Cơ sở về lý thuyết đồ thị

1.1 Đồ thị và một số khái niệm cơ bản

1.1.1 Đồ thị vô hướng

Định nghĩa 1.1.1. Đồ thị G là một đối tượng được cấu thành từ 3 thành phần: tập đỉnh $V(G)$, tập cạnh $E(G)$ và mối liên hệ giữa 2 đỉnh trong $V(G)$ (không nhất thiết phải phân biệt) thông qua một cạnh trong $E(G)$.

Định nghĩa 1.1.2. **Khuyên** là một cạnh nối từ một đỉnh vào chính nó. **Đa cạnh** là những cạnh có cùng một cặp đỉnh (có định hướng nếu trong đồ thị có hướng). Một đơn đồ thị là đồ thị không có khuyên và không có đa cạnh.

Ví dụ 1.1.1. Cho

$$G = (V(G), E(G), \psi_G)$$

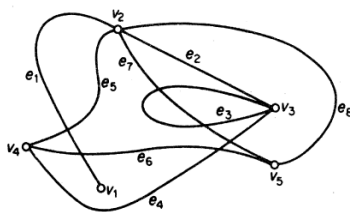
trong đó

$$\begin{aligned} V(G) &= \{v_1, v_2, v_3, v_4, v_5\} \\ E(G) &= \{e_1, e_2, e_3, e_4, e_5, e_6, e_7, e_8\} \end{aligned}$$

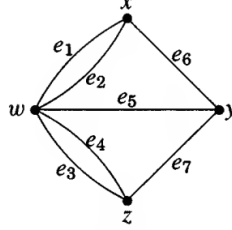
và ψ_G được xác định bởi

$$\psi_G(e_1) = v_1v_2, \psi_G(e_2) = v_2v_3, \psi_G(e_3) = v_3v_3, \psi_G(e_4) = v_3v_4$$

$$\psi_G(e_5) = v_2v_4, \psi_G(e_6) = v_4v_5, \psi_G(e_7) = v_2v_5, \psi_G(e_8) = v_2v_5$$



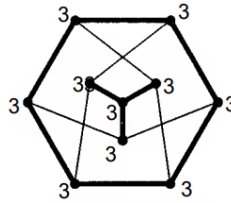
Định nghĩa 1.1.3. Xét đồ thị $G = (V, E)$ và cạnh $e = (u, v) \in E(G)$. Ta nói rằng hai đỉnh u và v là **kề nhau** và cạnh e này **liên thuộc** với hai đỉnh u và v .



Như trong đồ thị trên, ta nói rằng w và y là kề nhau và e_5 là liên thuộc với w, y .

Định nghĩa 1.1.4. Với một đỉnh $u \in V(G)$, ta ký hiệu **bậc của** u là $d_G(u)$ và bằng số cạnh liên thuộc với u , ngoại trừ khuyên do được tính hai lần. **Bậc lớn nhất** của đồ thị G là $\Delta(G)$, **bậc nhỏ nhất** là $\delta(G)$. Một đồ thị được gọi là **k -chính quy** nếu $\Delta(G) = \delta(G) = k$.

Định nghĩa 1.1.5. **Bậc của đồ thị** G , ký hiệu $n(G)$ là tổng số đỉnh trong đồ thị G . **Kích thước của đồ thị** G , ký hiệu $e(G)$ là tổng số cạnh trong G .



Hình 1.1: Đồ thị Petersen

Đồ thị Petersen trên là đồ thị 3-chính quy do $\Delta(G) = \delta(G) = 3$ và bậc của G bằng 10, kích thước của nó bằng 15.

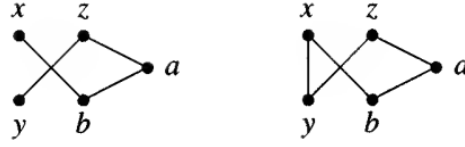
Định lý 1.1.1 (Công thức tổng bậc). Cho đơn đồ thị G , ta có

$$\sum_{v \in V(G)} d(v) = 2e(G)$$

Chứng minh. Do mỗi cạnh của đồ thị G đều liên thuộc với 2 đỉnh phân biệt nên phép lấy tổng các bậc sẽ đếm mỗi cạnh 2 lần. ■

Hệ quả 1.1.1.1. Mọi đồ thị đều có một số chẵn các đỉnh bậc lẻ. Không tồn tại một đồ thị với số lẻ các đỉnh là chính quy với bậc của đỉnh là lẻ.

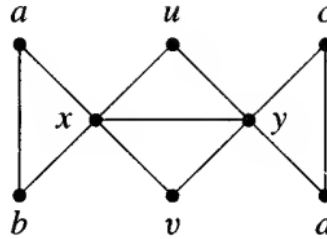
Định nghĩa 1.1.6. Một **path** là một đồ thị đơn với các đỉnh được sắp xếp sao cho 2 đỉnh là kề khi và chỉ khi chúng là liên tiếp nhau trong một danh sách nào đó. Chu trình là một path với đỉnh đầu và đỉnh cuối trùng nhau.



Hình 1.2: Path và Cycle

Định nghĩa 1.1.7. Cho đồ thị G . **Walk** là một danh sách các đỉnh và cạnh $v_0, e_1, v_1, \dots, e_k, v_k$ sao cho, với $1 \leq i \leq k$, thì cạnh e_i liên thuộc với v_{i-1} và v_i . **Trail** là một walk nhưng không có sự lặp lại về cạnh. Khi ta nói u, v -walk hay u, v -trail tức là u, v là 2 đỉnh đầu cuối của nó. Một u, v -**path** là một đường đi từ u đến v mà các đỉnh bậc 1 của nó là u và v . Các đỉnh còn lại ngoài u và v được gọi là các **đỉnh trong**.

Độ dài của một walk, trail, path, cycle là bằng số cạnh của nó. Ta gọi walk, trail là đóng (closed) nếu đỉnh đầu và cuối của nó trùng nhau.

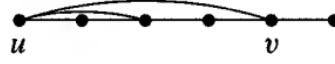


Ta minh họa các khái niệm trong hình trên. Dễ thấy một closed walk với độ dài 12 là danh sách $(a, x, a, x, u, y, c, d, y, v, x, b, a)$. Nếu bỏ đi 2 đỉnh đầu trong danh sách ta sẽ được một closed trail với độ dài 10.

Đồ thị trên có 5 chu trình: (a, b, x) , (c, y, d) , (u, x, y) , (x, y, v) , (u, x, v, y) . Danh sách (u, y, c, d, y, x, v) tạo thành một u, v -trail chứa các cạnh của một u, v -path u, y, x, v nhưng không chứa u, v -path u, y, v .

Định lý 1.1.2. Cho đồ thị G với $\Delta(G) = 2$ thì G chứa một chu trình.

Chứng minh. Giả sử P là path dài nhất có thể trong G và u là một đỉnh đầu cuối của P . Do P không thể được kéo dài thêm nữa nên mọi đỉnh kề với u phải là một đỉnh trong P . Mặt khác, u có số bậc nhỏ nhất bằng 2 nên tồn tại 1 đỉnh v kề với u thông qua một cạnh không thuộc P . Cạnh này sẽ tạo thành một chu trình trong đồ thị. ■

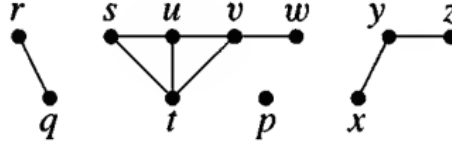


Định nghĩa 1.1.8. Một **đồ thị con** của đồ thị G là một đồ thị H với $V(H) \subseteq V(G)$ và $E(H) \subseteq E(G)$. Khi đó, ta có thể ký hiệu $H \subset G$.

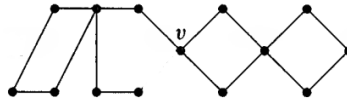
Đồ thị G được gọi là **liên thông** nếu mọi cặp đỉnh trong G đều thuộc một path; ngược lại, G là **không liên thông**.

Định nghĩa 1.1.9. Các **thành phần liên thông** của đồ thị G là các đồ thị con liên thông lớn nhất có thể của nó. Một đỉnh được gọi là **đỉnh cô lập** nếu bậc của nó bằng 0.

Đồ thị dưới đây có 4 thành phần liên thông và một đỉnh cô lập.



Định nghĩa 1.1.10. **Đỉnh cắt** hay **cạnh cắt** của một đồ thị là một đỉnh hoặc một cạnh mà nếu ta loại bỏ nó thì sẽ làm tăng số thành phần liên thông.



Hình 1.3: Đồ thị với đỉnh cắt v và cạnh cắt bên trái của v .

Định nghĩa 1.1.11. Một đồ thị được gọi là **Eulerian** nếu nó có một trail đóng chứa tất cả các cạnh. Một trail đóng còn được gọi là một **circuit** nếu chúng ta không muốn chỉ rõ đâu là đỉnh bắt đầu. Một **Eulerian circuit/Eulerian trail** trong một đồ thị là một circuit/trail chứa tất cả các cạnh.

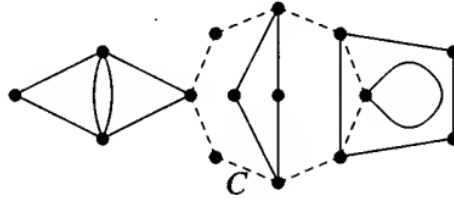
Định lý 1.1.3. Đồ thị G là Eulerian khi và chỉ khi nó là liên thông và các đỉnh của nó đều có bậc chẵn.

Chứng minh. Điều kiện cần: Hiển nhiên G phải liên thông. Giả sử G có một Eulerian circuit C . Dễ thấy muốn đi qua một đỉnh của C phải sử dụng 2 cạnh liên thuộc với nó do đó mọi đỉnh phải có bậc chẵn.

Điều kiện đủ: Với điều kiện liên thông và các đỉnh có số bậc chẵn. Chúng ta xây dựng một Eulerian circuit dựa vào quy nạp trên số cạnh m .

Bước cơ sở: $m = 0$, ta có một trail đóng bao gồm một đỉnh.

Bước quy nạp: $m > 0$. Với số bậc chẵn thì mọi đỉnh trong G có số bậc nhỏ nhất bằng 2. Theo định lý 1.1.2 thì G chứa một chu trình C . Cho $G' = G - E(C)$, hiển nhiên các thành phần của G' cũng là các đồ thị với số bậc của đỉnh là chẵn. Do các thành phần của G' cũng thỏa mãn điều kiện quy nạp nhưng với số cạnh ít hơn m nên ta áp dụng giả thuyết quy nạp dẫn đến mọi thành phần của G' đều có một Eulerian circuit. Để xây dựng được Eulerian circuit của G chúng ta duyệt qua C , nhưng khi đi đến đỉnh thuộc một thành phần của G' (lần đầu tiên) thì ta sẽ duyệt theo Eulerian circuit của thành phần đó rồi lại duyệt tiếp C . Circuit của G sẽ kết thúc tại đỉnh đầu tiên chúng ta bắt đầu duyệt C . ■



Hình 1.4: Eulerian circuit

1.1.2 Phép đẳng cấu trong đồ thị

Định nghĩa 1.1.12. Một phép **đẳng cấu** từ một đồ thị đơn G đến một đồ thị H là một song ánh $f : V(G) \rightarrow V(H)$ sao cho với mỗi $uv \in E(G)$ khi và chỉ khi $f(u)f(v) \in E(H)$. Lúc đó, ta nói rằng " G là đẳng cấu với H ", ký hiệu $G \cong H$.

Ví dụ 1.1.2. Cho G và H là 2 path có 4 đỉnh. Xác định hàm số $f : V(G) \rightarrow V(H)$ với $f(w) = a$, $f(x) = d$, $f(y) = b$, $f(z) = c$. Để chỉ ra f có là một đẳng cấu, ta có thể kiểm tra xem f có bảo toàn sự hiện diện của một cạnh qua hai đỉnh bất kì trong G . Chú ý rằng, nếu ta thực hiện vài biến đổi trên $A(G)$ bằng việc thay đổi thứ tự trên dòng và cột thành w, y, z, x ta sẽ được $A(H)$. Như vậy, f là một đẳng cấu. Một phép đẳng cấu khác biến w, x, y, z thành c, b, d, a .

G

H

	w	x	y	z
w	0	1	0	0
x	1	0	1	0
y	0	1	0	1
z	0	0	1	0

	w	y	z	x
w	0	0	0	1
y	0	0	1	1
z	0	1	0	0
x	1	1	0	0

	a	b	c	d
a	0	0	0	1
b	0	0	1	1
c	0	1	0	0
d	1	1	0	0

Chú ý 1.1.1. *Tìm kiếm phép đẳng cấu.* Ta có thể thực hiện một phép giao hoán cho cả dòng và cột của $A(G)$. Nếu ma trận mới thu được giống với $A(H)$ thì phép giao hoán đó trả về một phép đẳng cấu.

1.1.3 Phân hoạch trên đồ thị

1.1.4 Một số dạng đặc biệt của đồ thị

1.1.5 Đồ thị có hướng

1.2 Các cách biểu diễn đồ thị

Có rất nhiều cấu trúc dữ liệu có thể được sử dụng để biểu diễn đồ thị. Mỗi cách đều có ưu và nhược điểm riêng. Phần này sẽ thảo luận về một số cách biểu diễn động và tĩnh của đồ thị được tham khảo trong [2].

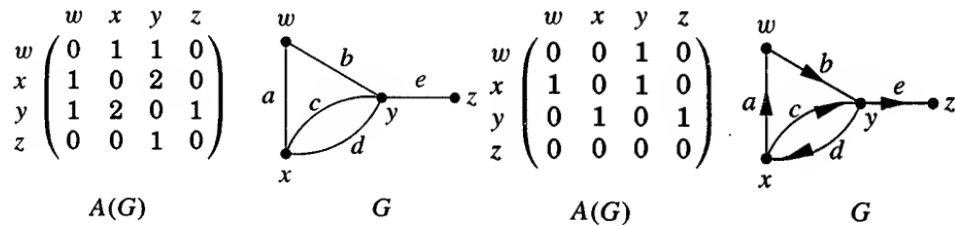
1.2.1 Biểu diễn tĩnh

Ma trận kề

Ma trận kề A của một đồ thị $G(V, E)$ là một ma trận cỡ $|V| \times |V|$ được xác định bởi:

$$A(i, j) = \begin{cases} 1, & \text{nếu } (i, j) \in E(G) \\ 0, & \text{nếu } (i, j) \notin E(G) \end{cases}$$

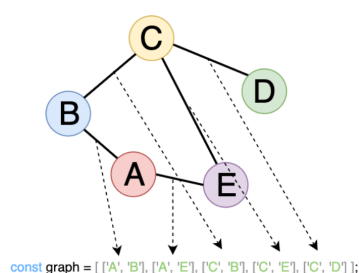
Cách biểu diễn trên có thể được dùng cho cả đồ thị vô hướng và có hướng. Tuy nhiên, nếu trong đồ thị vô hướng ma trận kề là đối xứng thì trong đồ thị có hướng điều này lại không đúng.



Hình 1.5: Ma trận kề của đồ thị vô hướng và có hướng (tham khảo trong [1])

Danh sách cạnh

Cách biểu diễn này chỉ đơn giản là một danh sách các cặp đỉnh không thứ tự (vô hướng) hoặc có thứ tự (có hướng).



Hình 1.6: Danh sách cạnh của một đồ thị vô hướng

Ma trận liên thuộc

1.2.2 Biểu diễn động - danh sách kề

1.3 Đồ thị dạng cây

1.3.1 Cây và các tính chất cơ bản

1.3.2 Cây bao trùm

1.4 Đường đi ngắn nhất

1.4.1 Thuật toán BFS

1.4.2 Thuật toán Dijkstra

1.4.3 Thuật toán Floyd-Warshall

1.4.4 Thuật toán Bellman-Ford

1.5 Mạng lưới và tính liên thông

Định nghĩa 1.5.1 (Mạng lưới). Một **mạng (network)** là một đồ thị có hướng với **sức chứa (capacity)** $c(e)$ không âm trên mỗi cạnh e và 2 điểm phân biệt: **nguồn s (source)** và **đích t (sink)**.

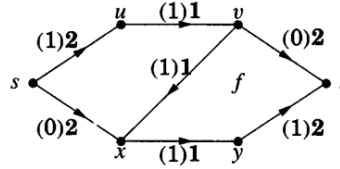
Thực tế cho thấy rất nhiều bài toán quan trọng đều có thể được mô hình hóa dưới dạng một mạng. Ví dụ như các tuyến đường giao thông, các đường ống nước, các đường dẫn dữ liệu trong máy tính hay các dòng điện trong một mạng lưới điện, ...

1.5.1 Luồng trong mạng

Định nghĩa 1.5.2 (Luồng). Một **luồng (flow)** f có thể hiểu như một hàm số gán một giá trị $f(e)$ cho mỗi cạnh e . Trong đó,

$$\begin{aligned} f^+(v) & \text{ là tổng luồng trên các cạnh rời khỏi } v \\ f^-(v) & \text{ là tổng luồng trên các cạnh đi vào } v \end{aligned}$$

Một luồng được gọi là **khả thi (feasible)** nếu nó thỏa mãn **điều kiện về sức chứa** $0 \leq f(e) \leq c(e)$ cho mỗi cạnh và **điều kiện bảo toàn** $f^+(v) = f^-(v)$ tại mỗi nút $v \notin \{s, t\}$.



Hình 1.7: Ví dụ về một luồng khả thi

(sức chứa được thể hiện bởi số in đậm và luồng giá trị là các số trong ngoặc)

Định nghĩa 1.5.3. Giá trị của một luồng ($\text{val}(f)$) là giá trị trên các cạnh chảy vào điểm đích $f^-(t) - f^+(t)$.

Ví dụ 1.5.1. Ở hình (1.7), ta có giá trị của luồng là $\text{val}(f) = f^-(t) - f^+(t) = 0 + 1 - 0 = 1$.

1.5.2 Luồng cực đại

Như ở trong định nghĩa (1.5.3), ta biết rằng mỗi luồng đều xác định một giá trị nào đó. Luồng cực đại là luồng (khả thi) làm cho giá trị đó đạt lớn nhất.

Định nghĩa 1.5.4 (Augmenting path). Cho f là một luồng khả thi trên một mạng N , một **đường tăng luồng** là một path từ điểm nguồn tới điểm đích P trong đồ thị G (ẩn dưới mạng N) sao cho với mỗi $e \in E(P)$:

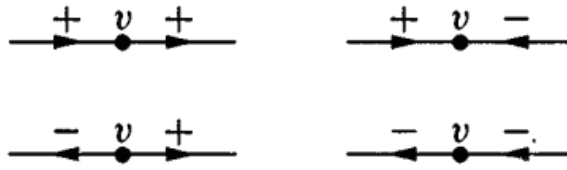
- nếu P đi cùng chiều với e , thì $f(e) < c(e)$.
- nếu P đi ngược chiều với e , thì $f(e) > 0$.

Đặt $\varepsilon(e) = f(e) - c(e)$ khi e cùng chiều với P , hoặc $\varepsilon(e) = f(e)$ nếu e ngược chiều với P . **Dung sai (tolerance)** trên P sẽ bằng $\min_{e \in E(P)} \varepsilon(e)$.

Định lý 1.5.1. Nếu P là một đường tăng luồng với dung sai z , thì việc thay đổi luồng bằng cách cộng z trên các cạnh cùng chiều với P và $-z$ trên các cạnh ngược chiều với P sẽ sinh ra một luồng mới f' (khả thi) với $\text{val}(f') = \text{val}(f) + z$.

Chứng minh. Do định nghĩa của dung sai luôn đảm bảo $0 \leq f'(e) \leq c(e)$ với mỗi cạnh e , nên điều kiện về sức chứa được thỏa mãn. Với điều kiện về bảo toàn trên luồng, ta chỉ cần kiểm các đỉnh trong v .

Các cạnh của P đi qua đỉnh trong của P nằm trong 1 trong 4 trường hợp sau.



Ở bất kì trường hợp nào thì sự thay đổi trên các luồng ra cũng đều tương tự với sự thay đổi trên các luồng vào. ■

Bài toán **luồng cực đại (max flow)** là một bài toán đối ngẫu với bài toán **lát cắt hẹp nhất (min cut)**. Mối liên hệ giữa hai bài toán này có thể được tìm hiểu thêm ở [1].

Dưới đây là một thuật toán do Ford-Fulkerson đề xuất, ý tưởng của nó là tìm kiếm một đường tăng luồng để làm tăng giá trị của luồng. Nếu không thể tìm được một đường nào như vậy, thì nó sẽ tìm ra một tập *min cut* ở trên với giá trị bằng với giá trị lớn nhất của luồng (**Max flow-min cut** [1]).

Thuật toán 1.5.1 (Thuật toán Ford-Fulkerson).

Đầu vào: Một luồng khả thi f trong một mạng.

Đầu ra: Một đường tăng luồng của f hoặc một tập với lực lượng bằng $\text{val}(f)$.

Ý tưởng: Tìm những nút có thể tới được từ nút s bằng những path có dung sai dương. Nếu có thể đi đến nút t thì nó tạo thành một đường tăng luồng. Trong quá trình tìm kiếm, R là tập các nút được đánh dấu *Reached*, và S là tập con của R được đánh dấu *Searched*.

Khởi tạo: $R = \{s\}$, $S = \emptyset$.

Vòng lặp: Chọn $v \in R - S$.

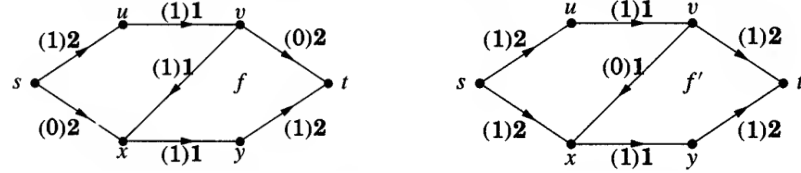
Với mỗi cạnh vw có $f(vw) < c(vw)$ và $w \notin R$ thì ta thêm w vào R .

Với mỗi cạnh uv có $f(uv) > 0$ và $u \notin R$ thì thêm u vào R .

Đánh nhãn mọi đỉnh thêm vào R là *Reached*, và ghi lại đỉnh v là đỉnh xuất phát. Sau khi xem xét tất cả cạnh liên thuộc với v , ta thêm v vào S .

Nếu điểm đích t là nằm trong R , thì ta truy vết đường mà đi đến t để trả về một đường tăng luồng f và kết thúc. Nếu $R = S$, thì trả về một tập *min cut*. Nếu không có cái nào thỏa mãn, thì ta lặp lại thuật toán.

Ví dụ 1.5.2. Chúng ta sẽ áp dụng thuật toán Ford-Fulkerson với ví dụ trong hình (1.7).



Đầu tiên, tìm kiếm từ s các sức chứa còn dư tới u và tới x , đánh nhãn chúng là *reached*. Khi đó, $u, v \in R - S$. Do không còn sức chứa trên cạnh uv và xy , nên ta không thể dùng các cạnh này. Tuy nhiên, trên cạnh vx có luồng giá trị bằng $1 > 0$ nên ta đánh dấu v từ nút x . Do đó, v là phần tử duy nhất trong $R - S$ và tìm kiếm từ đỉnh này giúp ta tới được t . Truy vết lại các đỉnh trước, ta được một đường tăng luồng s, x, v, t .

Dung sai trên s, x, v, t là 1, nên luồng mới có giá trị tăng lên 1 ($= 2$) được thể hiện ở hình bên phải. Nếu chạy thuật toán lần nữa, ta sẽ thấy rằng $R = S = \{s, u, x\}$ và ta dừng thuật toán tại đây với kết luận rằng luồng đã đạt cực đại với giá trị bằng 2.

Chương 2

Trực quan hóa đồ thị với NetworkX API

2.1 Giới thiệu về NetworkX API

NetworkX là một python package được thiết kế để tạo ra, thao tác và nghiên cứu về cấu trúc của dữ liệu dạng đồ thị. Đây là một công cụ hỗ trợ mạnh mẽ để giúp những người nghiên cứu về lý thuyết đồ thị dễ dàng hiện thực hóa các ý tưởng của họ.

Để làm quen với networkX, chúng em sẽ sử dụng bộ dữ liệu do Konekt Seventh Graders cung cấp. Đây là dữ liệu dạng đồ thị có hướng chứa các đánh giá lẫn nhau của 29 học sinh khối 7. Các đỉnh đại diện cho các học sinh. Các cạnh có hướng và được đánh trọng số (1 đến 3) thể hiện rằng một học sinh đánh giá một học sinh khác như thế nào.

2.1.1 Mô hình dữ liệu

Trong NetworkX, dữ liệu dạng đồ thị được lưu trữ dưới một cấu trúc kiểu từ điển **Graph**.

```
G = nx.Graph() # or nx.DiGraph()
```

Các đỉnh có thể được truy xuất thông qua thuộc tính **nodes** của **G**. Tương tự với các cạnh được truy xuất thông qua thuộc tính **edges**, tuy nhiên ta cần truyền vào 2 đỉnh phân biệt.

```
# accessing to node 1
G.nodes[node1]
# accessing to edge between node 1 and node 2
G.edges[node1, node2]
```


Do **Graph** được cài đặt dưới dạng một từ điển, nên bất kỳ một **hashable object** nào cũng có thể là một đỉnh (**string** hoặc **tuples**, nhưng không phải **list** hay **sets**)

Sau khi khởi tạo **Graph**, ta cần nạp dữ liệu vào trong đối tượng này bằng phương thức sau:

```
G.add_edge(u, v, **attr)
# u: first node
# v: second node
# **attr: attributes of edge
```

Đối với bộ dữ liệu *seventh grades* như đã nói ở trên, chúng em đã viết một hàm để đưa dữ liệu vào trong đối tượng **Graph**.

```
G = load_data.load_seventh_grader_network()
```

2.1.2 Các thống kê cơ bản của một đồ thị

Khi đã xây dựng được dữ liệu đồ thị, một trong những điều tiên ta cần làm là kiểm tra các thống kê cơ bản của nó, chẳng hạn như: số lượng các đỉnh, số lượng các cạnh, ...

Truy vấn loại đồ thị

Trong trường hợp ta không rõ đồ thị chúng ta đang làm việc là vô hướng hay có hướng thì có thể sử dụng câu lệnh sau để kiểm tra

```
type(G)
# return: networkx.classes.digraph.DiGraph
```

Truy vấn thông tin trên đỉnh

Để truy vấn tập hợp đỉnh ta dùng câu lệnh sau

```
G.nodes()
# return: NodeView((1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15,
16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29))
```

Câu lệnh trên trả về một *"view"* của các đỉnh. Ngoài ra, kiểu dữ liệu trả về của câu lệnh trên là kiểu **iterable** nên ta có thể dùng câu lệnh

```
len(G.nodes())
# return: 29
```

để đếm số đỉnh trong đồ thị.

Nếu trên đồ thị có gắn *metadata*, thì ta cũng có thể xem được chúng với câu lệnh

```
list(G.nodes(data=True))[0:5] # show 5 nodes
# return: [(1, {'gender': 'male'}),
#         (2, {'gender': 'male'}),
#         (3, {'gender': 'male'}),
#         (4, {'gender': 'male'}),
#         (5, {'gender': 'male'})]
```

Câu lệnh trên trả về kiểu **NodeDataView** - một dạng của từ điển. Dựa vào tính chất của từ điển, ta có dòng duyệt qua từng cặp key - value của nó. Sau đây là hàm đếm số học sinh nam và nữ.

```
def node_metadata(G):
    from collections import Counter

    mf_counts = Counter([d["gender"] for n, d in G.nodes(data=True)])
    return mf_counts

print(node_metadata(G))
# return: Counter({'female': 17, 'male': 12})
```

Truy vấn thông tin trên cạnh

Câu lệnh sau sẽ giúp chúng ta lấy được một danh sách các cạnh trên một đồ thị

```
list(G.edges())[0:5] # show 5 edges
# return: [(1, 2), (1, 3), (1, 4), (1, 5), (1, 6)]
```

Tương tự như trên đỉnh, phương thức **G.edges()** trả về một **EdgeView** mà có thể duyệt qua được. Do vậy, ta có thể đếm được số cạnh bằng câu lệnh sau

```
len(G.edges())
# return: 376
```

Do trên cạnh cũng có thể tồn tại *metadata* nên bằng cách truyền tham số **data=True**, ta sẽ thu được một **EdgeDataView**

```
list(G.edges(data=True))[0:5]
# return: [(1, 2, {'count': 1}),
#         (1, 3, {'count': 1}),
#         (1, 4, {'count': 2}),
#         (1, 5, {'count': 2}),
```

```
# (1, 6, {'count': 3})]
```

Thêm vào đó, nếu muốn truy xuất một cạnh bất kì ta có thể làm như sau

```
G.edges[15, 10] # edge between vertex 15 and 10
# return: {'count': 2}

G.edges[15, 16] # edge between vertex 15 and 16
# get an error because this edge doesn't exist.
```

Do **EdgeDataView** có dạng một từ điển nên ta có thể duyệt qua các cặp key-value của nó. Đoạn code sau đây thống kê số lần nhiều nhất mà một học sinh đánh giá học sinh khác

```
def edge_metadata(G):
    counts = [d["count"] for n1, n2, d in G.edges(data=True)]
    max_count = max(counts)
    return max_count

print(edge_metadata(G))
# return: 3
```

2.1.3 Thao tác trên đồ thị

Thêm đỉnh

```
G.add_node(node, node_data1=some_value, node_data2=some_value)
```

Giả sử trong bộ dữ liệu bị thiếu mất 2 học sinh (số 30 và số 31). Ta cần thêm vào bộ dữ liệu số 30 (nam) và số 31 (nữ).

```
def add_students(G):
    G = G.copy()
    G.add_node(30, gender="male")
    G.add_node(31, gender="female")
    return G
```

Thêm cạnh

```
G.add_edge(node1, node2, edge_data1=some_value, edge_data2=some_value)
```

Giữa 2 học sinh 30 và 31 ta vừa thêm, giả sử họ rất thân với nhau, nên ta thêm một cạnh giữa node 30 và node 31.

```
def add_student_rating(G):  
    G = G.copy()  
    G.add_edge(30, 31, count=3)  
    G.add_edge(31, 30, count=3)  
    return G
```

Khai phá dữ liệu

Giả sử trong số các học sinh trong bộ dữ liệu, có bạn rất yêu mến một bạn học sinh khác nhưng lại không được đáp lại. Nhiệm vụ ở đây là tìm ra các học sinh có một tình bạn "đơn phương" như vậy

```
def unrequitted_friendships(G):  
    losers = []  
    for n1, n2 in G.edges():  
        if not G.has_edge(n1, n2):  
            losers.append((n1, n2))  
    return losers
```

Tuy nhiên, có thể thấy rằng trong mạng lưới học sinh này không có bất kì ai như vậy.

2.2 Trục quan hóa đồ thị

Ai đó đã từng nói:

“A picture is worth a thousand words”

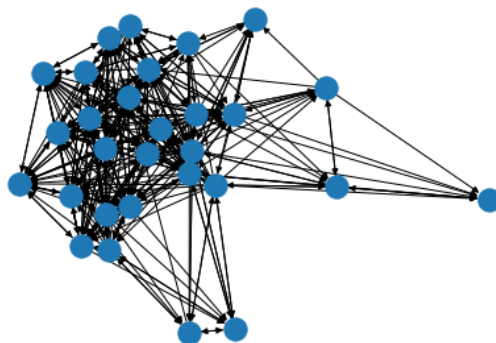
Đặc biệt trong lĩnh vực lý thuyết đồ thị, câu nói trên càng được khẳng định tính đúng đắn của nó. Việc sử dụng các hình ảnh hợp lý, có thể giúp tìm ra cấu trúc ẩn sau đồ thị, cách mà các thành phần trong đồ thị giao tiếp với nhau hay đôi khi là dễ dàng chứng minh một mệnh đề hoặc định lý nào đó.

Vẫn sử dụng bộ dữ liệu *seventh graders* ở phần (2.1). Dưới đây là một số các cách để biểu diễn mạng lưới này lên máy tính thông qua networkX.

2.2.1 Hairball

Biểu đồ thể hiện mối liên kết giữa các đỉnh qua cạnh (có hướng hoặc vô hướng) là một trong những biểu đồ phổ biến nhất. Các đỉnh thường được biểu diễn dưới dạng các hình tròn và với cạnh thì là đoạn thẳng (có hướng hoặc vô hướng) nối giữa 2 đỉnh có liên kết với nhau.

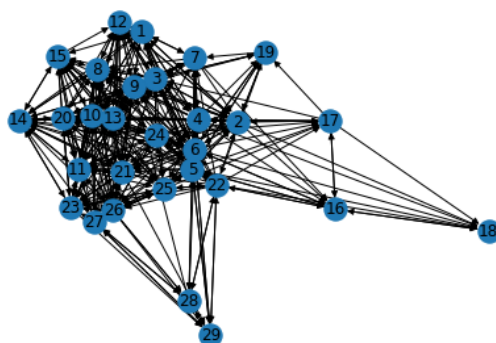
```
nx.draw(G)
```



NetworkX vẽ đồ thị dựa trên cơ chế chọn một đỉnh đầu tiên và vẽ các đỉnh, liên kết còn lại xung quanh đó. Do việc chọn đỉnh là ngẫu nhiên, nên ta thường sẽ không nhận được cùng một hình ở mỗi lần chạy. Tuy nhiên, có thể thấy rằng các đỉnh có xu hướng liên kết mạnh hơn thì sẽ được đặt gần nhau và tạo thành một cụm.

Ngoài ra, nếu đồ thị và các nhãn đủ nhỏ, ta có thể dùng câu lệnh sau để hiển thị một số thông tin lên đồ thị.

```
nx.draw(G, with_labels=True)
```



Tuy nhiên, với một đồ thị cỡ lớn thì cách biểu diễn này không thực sự là một sự lựa chọn tốt.

2.2.2 Matrix plot

Như trong chương 1 đã thảo luận, một trong những cách mà máy tính định nghĩa được một đồ thị là qua dạng ma trận.

```
nv.matrix(G, group_by="gender", node_color_by="gender")

from nxviz import annotate
annotate.matrix_group(G, group_by="gender")
```



Đối với bộ dữ liệu này, các đỉnh nằm trên các trục Ox, Oy là các học sinh và được nhóm theo giới tính bởi 2 màu (xanh là nữ, vàng là nam). Các hình tròn màu xám biểu thị giữa 2 học sinh có mối liên hệ với nhau (thông qua *đánh giá*). Ngoài ra, từ ma trận trên có thể thấy rằng:

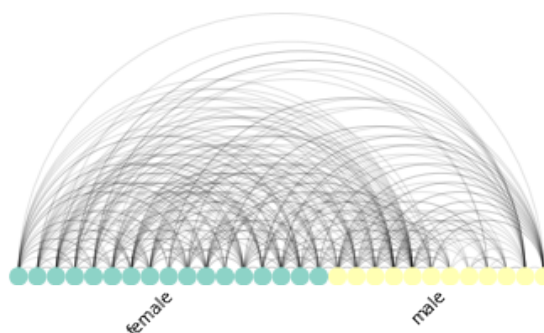
- không có học sinh nào tự bầu chọn cho chính mình (không có chấm nào trên đường chéo).
- ma trận trên không đối xứng do là đây là đồ thị có hướng

Cách biểu diễn trên cũng có một hạn chế là ta không thể biết được mức độ đánh giá của một học sinh này đối với một học sinh khác là như thế nào.

2.2.3 Arc plot

Arc plot là một dạng biểu diễn đồ thị đặc biệt. Trong đó, các đỉnh được biểu diễn bằng các thực thể hình tròn và các cạnh nối với nhau thể hiện mối quan hệ giữa 2 thực thể. Các đỉnh (hay thực thể) được sắp xếp trên một trục duy nhất và các liên kết thể hiện bởi các cạnh dạng vòng cung.

```
nv.arc(G, node_color_by="gender", group_by="gender")
annotate.arc_group(G, group_by="gender")
```

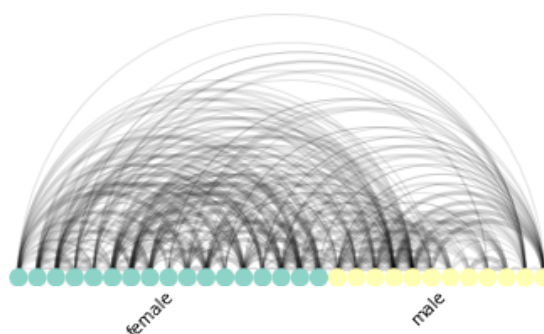


Kiểu biểu diễn arc plot có 2 điểm mạnh sau

- Nếu các đỉnh được sắp xếp tốt, arc plot có thể chỉ rõ các cụm hoặc các cầu.
- Việc gán nhãn các đỉnh là dễ dàng hơn so với những kiểu biểu diễn thông thường như hairball.

Một biến thể khác của arc plot là ta có thể làm cho các cạnh của nó lớn hơn nếu trong đồ thị là có trọng số.

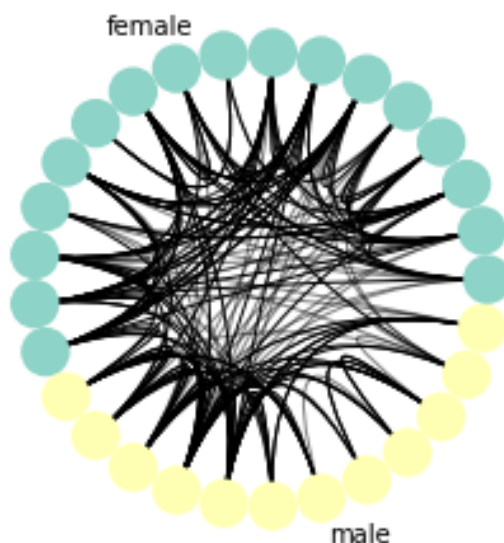
```
nv.arc(G, node_color_by="gender", group_by="gender", edge_lw_by="count")
annotate.arc_group(G, group_by="gender")
```



2.2.4 Circos plot

Circos là một sự cải tiến từ arc plot. Trong kiểu biểu diễn này, các đỉnh tạo với nhau thành một vòng tròn bao quanh các liên kết giữa các đỉnh.

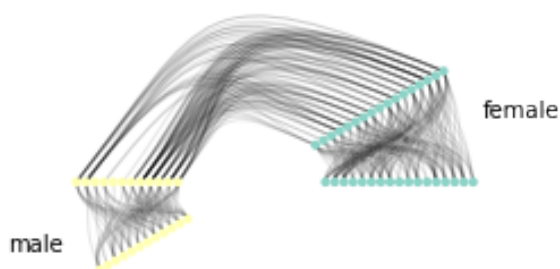
```
nv.circos(G, group_by="gender", node_color_by="gender",
edge_alpha_by="count")
annotate.circos_group(G, group_by="gender")
```



Cách biểu diễn dạng circos được sử dụng khá phổ biến và rộng rãi do nó đẹp và nó là lý tưởng để nhận biết mối quan hệ giữa các thực thể và cấu trúc ẩn sau của chúng.

2.2.5 Hive plot

```
nv.hive(G, group_by="gender", node_color_by="gender")
annotate.hive_group(G, group_by="gender")
```



Như trên hình có thể thấy, hive plot nhóm các đỉnh về 2 hoặc 3 trục đồng tâm. Trong trường hợp trên, các học sinh nam được đặt trên một trục và các bạn nữ ở trục còn lại.

Trước tiên, các cạnh được xây dựng dựa trên liên kết giữa các học sinh của 2 nhóm khác nhau. Sau đó, tiếp tục xây dựng các cạnh bên trong 1 nhóm bằng cách nhân bản các trục nam và nữ. Như trên hình, ta thấy *male* có 2 nhóm với cùng số thực thể và *female* cũng vậy.

Chương 3

Ứng dụng trong mạng lưới giao thông của Uber

Uber Technologies, Inc. (Uber) là một công ty cung cấp dịch vụ vận chuyển của Mỹ và có mặt trên toàn cầu. Như rất nhiều công ty lớn khác, để phục vụ cho mục đích cộng đồng, Uber cung cấp một lượng lớn dữ liệu vận chuyển. Trong bài báo cáo này, chúng em sử dụng tập dữ liệu do Uber cung cấp trên thành phố San Francisco của Mỹ trong quý 4 năm 2017 để làm quen với mô hình đồ thị trong thực tế.

Chuẩn bị dữ liệu:

1. Dữ liệu do Uber cung cấp có thể được tìm thấy tại "*Uber Movement*".
2. Tìm thành phố San Francisco, chọn travel times và nhấn download.
3. Nếu download thành công, tên file sẽ có dạng `san_francisco-censustracts-2017-4-All-MonthlyAggregate.csv`
4. Ngoài ra, ta cũng cần download file Geo Boundaries từ trang web đó.

Mô tả dữ liệu: Nhìn chung, dữ liệu chúng ta có được chia làm 2 file: *Travel times* và *Geo boundaries*.

- File thứ nhất *Travel times* gồm các thống kê của Uber về thời gian di chuyển giữa các cặp vị trí (points) trong khu vực San Francisco. Các vị trí trên bản đồ được biểu diễn bởi các số ID duy nhất.
- File thứ hai *Geo boundaries* cho biết mối liên hệ giữa các ID và khu vực trên San Francisco. File này chứa kinh độ và vĩ độ của các đỉnh trên mỗi khu vực (hình đa giác). Ví dụ, nếu một khu vực được biểu diễn bởi một đa giác gồm 5 đỉnh thì dữ liệu của ta sẽ có dạng một ma trận cỡ 5×2 với các hàng là chỉ số của các đỉnh và 2 cột là vị trí vĩ độ và kinh độ.

3.1 Xây dựng mô hình đồ thị

Sử dụng thư viện *pandas*, đọc 2 file mà ta đã tải về vào trong 2 dataframe *data* và *data_json*.

```
data = pd.read_csv(
    datasets /
    "uber/san_francisco-censustracts-2017-4-All-MonthlyAggregate.csv",
    header=None,
    sep=","
)
json_data = pd.read_json(datasets /
    "uber/san_francisco-censustracts.json")
```

Mặc dù có nhiều cách di chuyển từ địa điểm *A* đến *B* và từ *B* về *A*, nhưng do chỉ quan tâm đến thời gian trung bình nên chúng em sẽ áp dụng mô hình đồ thị vô hướng.

```
G = nx.Graph()
```

Các đỉnh trên đồ thị tương ứng là các địa điểm trên bản đồ. Ngoài ra, trên mỗi đỉnh còn có 2 thuộc tính.

1. Display name: địa chỉ con đường qua khu vực đó.
2. Location: trung bình tọa độ các đỉnh của đa giác (vector 2 chiều).

```
node_dict = {}
for node_info in json_data["features"]:
    node_id = int(node_info["properties"]["MOVEMENT_ID"])
    if node_id in node_dict:
        pass
    else:
        dis_name = node_info["properties"]["DISPLAY_NAME"]
        location =
            np.mean(np.array(node_info["geometry"]["coordinates"][0][0]),
                axis=0)
        node_dict[node_id] = {"DISPLAY_NAME": dis_name, "Location":
            location}
        G.add_node(node_id, name=dis_name, Location=location)
```

Chúng em xử lý các cạnh bằng cách xét 2 điểm kề nhau, tính tổng trọng số trên mọi con đường qua 2 điểm đó và chia cho số đường.

```
data_row = data.shape[0]
edge_dic = {}
```

```

for idx in range(data_row):
    if int(data['month'][idx]) != 12:
        continue
    edge_w = float(data["mean_travel_time"][idx])
    source_id = int(data["sourceid"][idx])
    destin_id = int(data['dstid'][idx])
    if (source_id, destin_id) in edge_dic:
        edge_dic[(source_id, destin_id)][0] += edge_w
        edge_dic[(source_id, destin_id)][1] += 1
    elif (destin_id, source_id) in edge_dic:
        edge_dic[(destin_id, source_id)][0] += edge_w
        edge_dic[(destin_id, source_id)][1] += 1
    else:
        edge_dic[(source_id, destin_id)] = [edge_w, 1]

# add edges to graph
for key, item in edge_dic.items():
    w = item[0]/item[1]
    G.add_edge(key[0], key[1], weight=w)

```

Đồ thị lúc này có dạng như sau:



Do chỉ quan tâm đến các thành phần liên thông, nên chúng em sẽ loại bỏ các điểm cô lập và chỉ giữ lại thành phần lớn nhất

```

Gcc = max((G.subgraph(c) for c in nx.connected_components(G)), key=len)

print("Number of nodes are:", nx.number_of_nodes(Gcc)) # 1898
print("Number of edges are:", nx.number_of_edges(Gcc)) # 321703

```

Sau khi đã tiền xử lý xong, ta có thể thấy **Gcc** là một đồ thị liên thông với 1898 đỉnh và 321703 cạnh. Đây cũng là đồ thị mà chúng em tập trung phân tích ở những phần sau.

3.2 Bài toán người du lịch

Do thuộc lớp bài toán NP-C nên việc giải quyết bài toán người du lịch theo những thuật toán truyền thống là rất khó khăn. Vậy nên, ở đây chúng em xin đề xuất thuật toán *1-approximation*^[4] thuộc vào lớp các phương pháp xấp xỉ để giải bài toán người du lịch.

Bước 1: Trước tiên, để có thể áp dụng thuật toán *1-approximation* ta cần kiểm tra giả thuyết bất đẳng thức tam giác trên **Gcc**. Việc kiểm tra tất cả các tam giác trong đồ thị là không cần thiết, vậy nên chúng em chỉ lấy mẫu khoảng 1000 tam giác. Nếu số lượng tam giác thỏa mãn đủ lớn, ta sẽ ngầm hiểu rằng kết quả sinh ra bởi thuật toán này là một xấp xỉ tốt.

```
import random

bound = max(node_dict.keys())
sample_num = 0
triangle_num = 0
used = set()

while sample_num < 1000:
    v1 = random.randint(1, bound)
    v2 = random.randint(1, bound)
    v3 = random.randint(1, bound)

    if (v1 != v2 != v3) \
        and (v1, v2) in edge_dict and (v2, v3) in edge_dict and (v3, v1)
        in edge_dict and \
        (v1, v2, v3) not in used:
        used.add((v1, v2, v3))
        sample_num += 1

        d12 = edge_dict[(v1, v2)][0] / edge_dict[(v1, v2)][1]
        d23 = edge_dict[(v2, v3)][0] / edge_dict[(v2, v3)][1]
        d31 = edge_dict[(v3, v1)][0] / edge_dict[(v3, v1)][1]

        if (d12 + d23 > d31) and (d23 + d31 > d12) and (d12 + d31 > d23):
            triangle_num += 1

print("The percentage of triangles in the graph is",
      triangle_num/sample_num * 100)
# The percentage of triangles in the graph is 93.2
```

Bước 2: Từ đồ thị **Gcc**, ta sẽ xây dựng một cây con bao trùm nhỏ nhất. Có một vài thuật toán để xây dựng cây con bao trùm nhỏ nhất, tuy nhiên chúng

em sẽ áp dụng thuật toán của Kruskal do tác giả của cuốn [1] đề xuất.

```
MST = nx.minimum_spanning_tree(Gcc, algorithm='kruskal')
```

Với **MST** đã xây dựng ở trên, ta ghi lại một số địa chỉ kết nối với nhau thông qua các cạnh trên cây con bao trùm nhỏ nhất.

```
for point in sorted(list(nx.edges(MST))[0:5]):
    print(node_dict[point[0]]["DISPLAY_NAME"], "----",
          node_dict[point[1]]["DISPLAY_NAME"])
# 400 Northumberland Avenue, Redwood Oaks, Redwood City ---- 1500 Oxford
# Street, Palm Park, Redwood City
# 400 Northumberland Avenue, Redwood Oaks, Redwood City ---- 100 Fifth
# Avenue, South Fair Oaks, Redwood City
# 18300 Sutter Boulevard, Morgan Hill ---- 17300 Lotus Way, Morgan Hill
# 18300 Sutter Boulevard, Morgan Hill ---- 1900 Alpet Drive, Morgan Hill
# 3200 Huntsman Drive, Rosemont Park, Sacramento ---- 8900 Cal Center
# Drive, Sacramento
```

Dựa vào google map, có thể thấy các kết quả trên khá đáng tin. Với 2 địa chỉ gần nhau thì sẽ có thời gian di chuyển nhỏ tương đương với trọng số trên cạnh cũng nhỏ. Đây chính xác là một trong những tính chất của cây con bao trùm nhỏ nhất.

Bước 3: Xây dựng một đồ thị đa cạnh từ **Gcc** bằng cách gấp đôi số cạnh trên mọi $uv \in E(\mathbf{Gcc})$

```
Multi_G = nx.MultiDiGraph()
Multi_G.add_nodes_from(MST)
path = []

for edge in MST.edges:
    w = MST.edges[edge]['weight']
    Multi_G.add_edge(edge[0], edge[1], weight=w)
    Multi_G.add_edge(edge[1], edge[0], weight=w)
```

Bước 4: Tìm một chu trình Eulerian trong **Multi_G**. Từ đó xây dựng phương án di chuyển cho bài toán người du lịch.

```
Euler_circle = nx.eulerian_circuit(multi_G)

travel_length = 0
idx = 0

for edge in Euler_circle:
    if idx == 0:
        start = edge[0]
    if edge[1] == start:
```

```

        print(edge[0])
        dst = edge[0]
        path.append(dst)
        print('-'*30)
        break

    travel_length += multi_G.edges[edge[0], edge[1], 0]['weight']
    path.append(edge[0])
    print(edge[0], "->", end=' ')
    idx += 1

print("Approximate TSP cost is", travel_length)
# return: Approximate TSP cost is 440276.24000000001

```

Bước 5: Đánh giá tính hiệu quả của thuật toán. Gọi ρ là hiệu suất về mặt chi phí của phương án xấp xỉ so với phương án tối ưu. Ta có:

$$\rho = \frac{\text{Chi phí của phương án xấp xỉ}}{\text{Chi phí của phương án tối ưu}}$$

Mặt khác,

$$\begin{aligned} \text{chi phí của phương án xấp xỉ} &\leq \text{chi phí trên chu trình Euler} \\ &\leq 2 \times \text{chi phí trên cây bao trùm nhỏ nhất} \\ &\leq 2 \times \text{chi phí trên phương án tối ưu nhất} \end{aligned}$$

Do vậy,

$$\rho \leq \frac{\text{chi phí trên chu trình Euler}}{\text{chi phí trên cây bao trùm nhỏ nhất}}$$

```

opt_length = 0

for edge in MST.edges:
    opt_length += MST.edges[edge]['weight']

print("The upper bound on the empirical performance of the approximation
      algorithm is", travel_length/opt_length)
# The upper bound on the empirical performance of the approximation
  algorithm is 1.5217849499512972

```

Bước 6: Trực quan hóa phương án di chuyển xấp xỉ tối ưu của bài toán người du lịch.

```

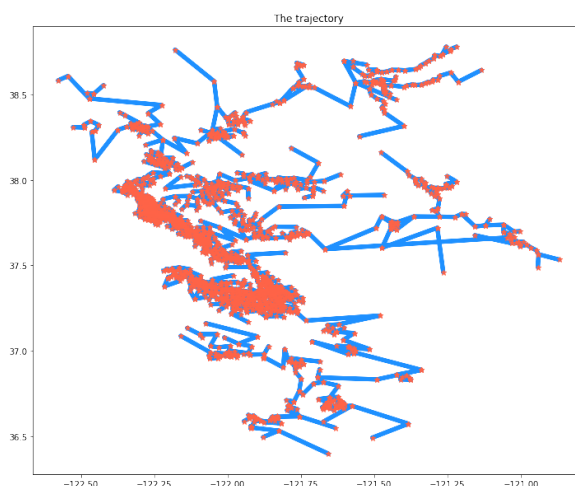
X = [] # latitude
Y = [] # longitude

for node in path:
    X.append(node_dict[node]['Location'][0])

```

```
Y.append(node_dict[node]['Location'][1])
```

```
plt.figure(figsize=(12, 10))
plt.plot(X, Y, linewidth=5, color='dodgerblue')
plt.plot(X, Y, '*', markersize=7, color='tomato')
plt.title('The trajectory')
```

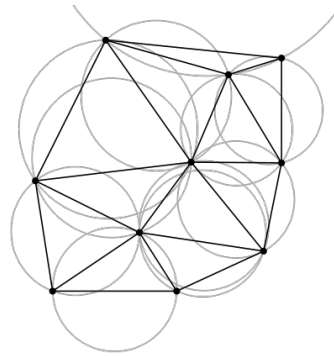


3.3 Khái quát hóa bản đồ giao thông

Khái quát hóa bản đồ là quá trình tạo nên một bản đồ dễ đọc từ một bộ dữ liệu địa lý chi tiết. Điều này đạt được bằng cách loại bỏ một số chi tiết và chỉ giữ lại những thứ quan trọng. Ví dụ như loại bỏ một số con đường (cạnh), khu vực; kết hợp một số khu vực lân cận; phóng đại các đối tượng địa lý; giảm kích thước các con đường, khu vực; dịch chuyển để đảm bảo khoảng cách giữa các đối tượng là phù hợp. Các phép toán trên có thể được tham khảo trong [5], và sâu hơn trong một mạng lưới giao thông thực sự [7].

Một trong những phương thức được sử dụng để khái quát hóa bản đồ là dùng *Delaunay Triangulation* (thuật toán về nó được thảo luận trong [6]).

Định nghĩa 3.3.1 (Tam giác phân Delaunay). Cho một tập P các điểm rời rạc, và một tập cạnh nối các điểm đó thành một tập các tam giác. Tập tam giác đó được gọi là tam giác phân Delaunay nếu không tồn tại bất kì điểm nào nằm trong 1 đường tròn ngoại tiếp tam giác bất kỳ.

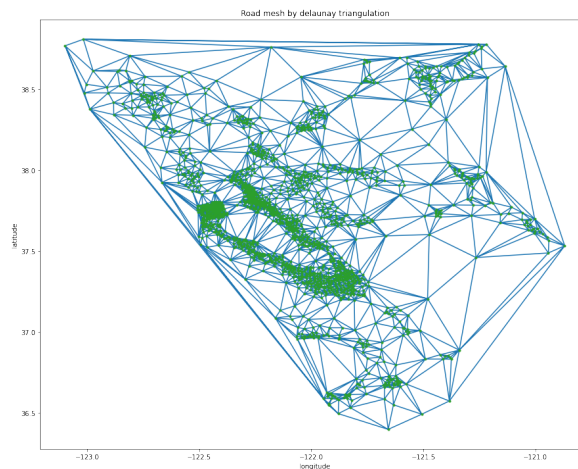


Do cài đặt thuật toán khá khó khăn nên chúng em sẽ sử dụng gói có sẵn trong python.

```
from scipy.spatial import Delaunay

nums_of_area = len(node_dict)
vertices = nx.nodes(Gcc)
mean_coordinates = np.array([elm['Location'] for elm in
                             list(node_dict.values())])
mean_coordinates_gcc = mean_coordinates[[x-1 for x in vertices], :]
triangulation = Delaunay(mean_coordinates_gcc)

plt.figure(figsize=(15, 12))
plt.triplot(mean_coordinates_gcc[:, 0], mean_coordinates_gcc[:, 1],
            triangulation.simplices)
plt.plot(mean_coordinates_gcc[:, 0], mean_coordinates_gcc[:, 1], '.')
plt.xlabel('longitude')
plt.ylabel('latitude')
plt.title('Road mesh by delaunay triangulation')
```



Gọi đồ thị sau khi áp dụng thuật toán tam giác phân Delaunay là G_{Δ} . Các đỉnh được thể hiện bởi các chấm màu xanh lá cây và các cạnh là những đoạn thẳng

màu xanh nước biển. Các tam giác trong G_Δ cũng chính là số lượng những tam giác thỏa mãn bất đẳng thức tam giác mà ta nhắc đến trong bài toán người du lịch.

```
G_delta = nx.Graph()

used_nodes = triangulation.simplices
edges_delta = set()

for row in range(used_nodes.shape[0]):
    for idx in range(3):
        node1 = list(vertices)[used_nodes[row][idx]]
        node2 = list(vertices)[used_nodes[row][(idx+1) % 3]]
        if (node1, node2) not in edges_delta and (node2, node1) not in edges_delta:
            edges_delta.add((node1, node2))
            G_delta.add_edge(node1, node2)

print("Number of nodes in G_delta is", nx.number_of_nodes(G_delta))
print("Number of edges in G_delta is", nx.number_of_edges(G_delta))
# Number of nodes in G_delta is 1898
# Number of edges in G_delta is 5680
```

Sau khi có được G_Δ , ở các phần sau khi thực hiện tính toán về luồng ta sẽ sử dụng nó thay cho G_{cc} .

3.4 Lưu lượng giao thông

Giả sử ta có các thông tin sau

- Mỗi một độ trong kinh độ và vĩ độ tương đương với 69 dặm Anh.
- Chiều dài của mọi chiếc xe là 5 mét ≈ 0.003 dặm.
- Mỗi xe có một khoảng cách an toàn tối thiểu là 2 giây đối với chiếc xe kề với nó.
- Mọi con đường đều có 2 làn.

Với giả thuyết rằng không có tắc đường, và xem xét luồng giao thông mà ta tính toán được như sức chứa cực đại trên con đường đó.

Câu hỏi: Làm thế nào để tính toán được lưu lượng xe trên một giờ trên mỗi con đường?

Dựa vào các giả thuyết trên, ta có được các hằng số sau

```

miles_per_degree = 69
car_length = 0.003
safe_time = 2/3600
lanes = 2

```

Ngoài ra, ta cũng khai báo thêm 2 biến để lưu lại kết quả tính được.

```

traffic_flow = {}
fake_flow = []

```

Sở dĩ, ta có biến `fake_flow` là do chúng em đang xét trên đồ thị G_Δ nên giữa 2 địa điểm có thể tồn tại một con đường không nằm trong tập dữ liệu.

Để tính toán lưu lượng giao thông trên mọi con đường trong G_Δ , ta cần duyệt qua mọi cạnh của nó. Nếu có một cạnh giữa 2 đỉnh u và v trong G_Δ , ta sẽ đặt `mean_time` là thời gian di chuyển trung bình giữa 2 điểm (đơn vị theo giây). Nếu không tồn tại một cạnh như vậy, đơn giản là ta đặt `mean_time` bằng một số rất lớn giữa 2 điểm đó. Sau đó, sử dụng các công thức sau để tính toán lưu lượng giao thông trên mỗi con đường.

- Khoảng cách giữa các đỉnh = $\sqrt{(u_x - v_x)^2 + (u_y - v_y)^2}$
- tốc độ của xe trên 1 đoạn đường = $\frac{\text{khoảng cách trên đoạn đường đó}}{\text{thời gian di chuyển trung bình trên đoạn đường đó}}$.
- khoảng cách giữa 2 xe = tốc độ của xe \times khoảng cách an toàn.
- Lưu lượng xe = $\frac{\text{tốc độ xe}}{(\text{chiều dài xe} + \text{khoảng cách giữa 2 xe})} \times \text{số lần}$.

Sau khi biết được lưu lượng xe cộ trên từng con đường trong G_Δ , ta thêm các giá trị đó vào các cạnh như là một thuộc tính trên cạnh đó (tương đương với sức chứa của con đường).

```

edges_delta = nx.edges(G_delta)
for edge in list(edges_delta):
    u = edge[0]
    v = edge[1]

    mean_time = 1e8 # threshold
    if (u, v) in edge_dict:
        mean_time = edge_dict[(u,v)][0] / edge_dict[(u,v)][1]
    elif (v, u) in edge_dict:
        mean_time = edge_dict[(v,u)][0] / edge_dict[(v,u)][1]

    miles = np.sqrt(
        np.square(abs(node_dict[u]['Location'][0] -
            node_dict[v]['Location'][0])) * miles_per_degree) +

```

```

        np.square(abs(node_dict[u]['Location'][1] -
                        node_dict[v]['Location'][1]) * miles_per_degree)
    )

    car_speed = miles / (mean_time/3600)
    safety_distance = car_speed * safe_time
    car_flow = (car_speed / (car_length + safety_distance)) * lanes

    G_delta.add_edge(u, v, capacity=car_flow)
    traffic_flow[(u,v)] = car_flow
    if mean_time == 1e8:
        fake_flow.append(car_flow)

```

Ta có thể kiểm tra một vài kết quả qua đoạn code dưới đây.

```

from itertools import islice

def take(n, iterable):
    return list(islice(iterable, n))

n_roads = take(5, traffic_flow.items())
for i in range(len(n_roads)):
    print(f"{n_roads[i][0]}: {n_roads[i][1]}")

# (1444, 2017): 0.5261317608739441
# (1444, 583): 1.3213564449322088
# (1444, 1231): 1.1921773439325514
# (1444, 793): 0.954980147325594
# (1444, 38): 3324.9516604562614

```

Ý nghĩa của kết quả trên, lấy ví dụ trên đoạn đường giữa điểm 1444 và điểm 38, tức là trung bình trong một giờ có khoảng 3325 xe lưu hành trên đoạn đường đó.

Ở 4 kết quả đầu tiên có thể thấy rằng số lượng xe lưu thông trên đường chỉ khoảng 1 đến 2 xe trên một giờ. Và như đã nói ở trên, chúng ta đang làm việc với G_Δ , nên có thể sẽ tồn tại những cung đường ảo. Ta nên đặt một giá trị giới hạn để tách biệt những con đường này.

```

fake_flow_threshold = max(fake_flow)
print("The fake flow are all under", fake_flow_threshold) # 3.5027891
print("Number of fake flows are:", len(fake_flow)) # 392

```

Vậy tất cả những con đường mà có lưu lượng giao thông dưới 3.5027891 đều là những con đường ảo và có tổng cộng tất cả 392 con đường như vậy.

3.5 Bài toán về luồng cực đại

Xem xét 2 địa điểm sau:

- Địa điểm nguồn: 100 Campus Drive, Stanford
- Địa điểm đích: 700 Meder Street, Santa Cruz

Tính số lượng xe cộ tối đa có thể đi trong một giờ từ Stanford tới UCSC. Đồng thời, tính số lượng các path phân biệt (theo cạnh) giữa 2 địa điểm trên.

Để trả lời câu hỏi thứ nhất, trước tiên chúng em sẽ xác định các nút đại diện cho Stanford và UCSC.

```
Stanford_ID = 0
UCSC_ID = 0

for key in node_dict:
    if node_dict[key]['DISPLAY_NAME'] == '100 Campus Drive, Stanford':
        Stanford_ID = key
        print("Stanford node is", Stanford_ID)
    elif node_dict[key]['DISPLAY_NAME'] == '700 Meder Street, Santa Cruz':
        UCSC_ID = key
        print("UCSC node is", UCSC_ID)

# Stanford node is 2607
# UCSC node is 1968
```

Tiếp theo, chúng em sẽ sử dụng một hàm có sẵn của *networkx* để tìm luồng cực đại.

```
flow_value, flow_dict = nx.maximum_flow(G_delta, UCSC_ID, Stanford_ID)
print("The maximum flow between Stanford and UCSC is", flow_value)
# The maximum flow between Stanford and UCSC is 14866.477294089982
```

Hàm `nx.maximum_flow` ở trên trả về một biến giá trị (giá trị của luồng cực đại) và một từ điển chứa giá trị của luồng cực đại đó gán cho mỗi cạnh.

Tương tự, để tìm các path phân biệt (theo cạnh), chúng em dùng hàm `nx.edge_disjoint_paths`.

```
edge_disjoint_paths = nx.edge_disjoint_paths(G_delta, UCSC_ID,
                                             Stanford_ID)

count = 0
for path in edge_disjoint_paths:
    count += 1
```

```

print(path)
print("Number of edge-disjoint path is", count)

# [1968, 2241, 1980, 2242, 744, 1869, 1363, 2240, 2607]
# [1968, 1424, 1431, 1980, 1763, 1762, 1209, 1733, 1725, 2607]
# [1968, 1431, 1989, 938, 744, 1737, 1736, 1726, 2607]
# [1968, 748, 2241, 1171, 1955, 1763, 1737, 1363, 2607]
# [1968, 1980, 1955, 2458, 1210, 1762, 1736, 2607]
# Number of edge-disjoint path is 5

```

Để dễ dàng hình dung và tưởng tượng, ta có thể phóng to bản đồ khu vực Stanford và UCSC.

```

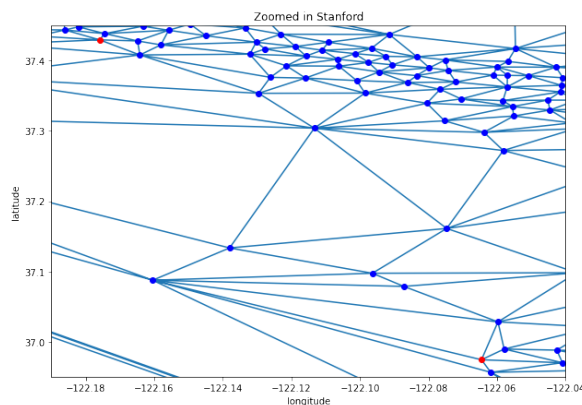
plt.figure(figsize=(10, 7))
plt.plot(Stanford_loc[0], Stanford_loc[1], 'o', color='red')
plt.plot(UCSC_loc[0], UCSC_loc[1], 'o', color='red')

for node in node_dict:
    if node != Stanford_ID and node != UCSC_ID:
        plt.plot(node_dict[node]['Location'][0],
                 node_dict[node]['Location'][1], 'o', color='blue')

plt.triplot(mean_coordinates_gcc[:, 0], mean_coordinates_gcc[:, 1],
            triangulation.simplices)
plt.xlim(-122.19, -122.04)
plt.ylim(36.95, 37.45)

plt.xlabel('longitude')
plt.ylabel('latitude')
plt.title('Zoomed in Stanford')
plt.show()

```



3.6 Cắt tĩa trên đồ thị

Kết luận

một vài kết luận

Tài liệu tham khảo

- [1] West, D. B. (2000). *Introduction to Graph Theory*. Prentice Hall. ISBN: 0130144002.
- [2] McHugh, J. A. (1990). *Algorithmic graph theory*. Prentice Hall. ISBN: 978-0-13-019092-5
- [3] Bondy, J. A., Murty, U. S. R. (1976). *Graph Theory with Applications*. New York: Elsevier.
- [4] Papadimitriou, C. H., & Steiglitz, K. (1998). *Combinatorial optimization: algorithms and complexity*. Courier Corporation.
- [5] Poorten, P. V. D., & Jones, C. B. (2002). *Characterisation and generalisation of cartographic lines using Delaunay triangulation*. International Journal of Geographical Information Science, 16(8), 773-794.
- [6] Shewchuk, J. R. (1997). *Delaunay refinement mesh generation*. Carnegie Mellon University.
- [7] Zhang, C., Li, Y., Xiang, L., Jiao, F., Wu, C., & Li, S. (2021). *Generating road networks for old downtown areas based on crowd-sourced vehicle trajectories*. Sensors, 21(1), 235.