

IDATT2101 Øving 5

Uvektede grafer

Trang Minh Duong

Deloppgave 1: Bredde-først søk

For å implementere en bredde-først søk implementerte en FIFO-kø med startnode s , merk $\text{dist}[s] = 0$. Når en node u tas ut av køen, besøkes alle nabover v som ikke er oppdaget, disse får $\text{dist}[v] = \text{dist}[u] + 1$ og $\text{parent}[v] = u$, og legges i køen. Siden køen er FIFO, prosesseres noder i ikke-sykende distanse, og første gang vi oppdager v har vi en korteste sti (i antall kanter).

```
1 BFS bfs(const vector<int>& offsets, const vector<int>& to,
2         int n, int s) {
3     BFS out{vector<int>(n, -1), vector<int>(n, -1)};
4
5     queue<int> q;
6     out.dist[s] = 0;
7     q.push(s);
8
9     while(!q.empty()) {
10        //Dequeue
11        int u = q.front();
12        q.pop();
13        out.order.push_back(u);
14
15        for(int x : adj[u]) {
16            if(out.dist[x] == INF) {
17                out.dist[x] = out.dist[u] + 1;
18                out.parent[x] = u;
19                q.push(x);
20            }
21        }
22    }
23    return out;
24 }
```

Listing 1: BFS og kø-implementasjon

Kompleksitet

Tidskompleksiteten er $O(N + M)$, ettersom hver node besøkes og hver kant behandles én gang. For en graf som lagres som `vector<vector<int>> adj`, er det slik at lagring av alle nabover: $O(M)$ (én oppføring per kant), vektorhoder for N noder: $O(N)$, og tillegg fra BFS: `dist[N]`, `parent[N]` og køen gir inntil $O(N)$. Totalt plasskompleksitet er da $O(N + M)$.

Lesing og konstruksjon av grafen er $O(N + M)$. Å lese inn antall noder N og kanter M fra fil, samt alle kantene (u, v) , krever $O(N + M)$ tid, lineært i størrelsen på input. Å bygge nabolistene ved å legge til hver kant én gang koster $O(M)$ tid, slik at total bygging er $O(N + M)$.

Begrensninger

Det er fortsatt mange begrensninger i besvarelsen, som:

Filformat: Første linje i filen skal være $N \ M$ (antall noder og kanter). Deretter følger nøyaktig M linjer med kanter $u \ v$. Programmet oppretter `adj(N)` og bruker `u` som indeks. Hvis `u` eller `v` er utenfor $[0, N-1]$, vil det oppstå indeksfeil. Besvarelsen kommer fra antagelsen at filen beskriver en rettet graf. For en urettet graf må begge retninger oppgis, eller legges til i koden

Uoppnåelige noder: Noder som ikke kan nås fra startnoden får `dist = -1` og `parent = -1`. I koden starter alle distanser som `INF` (uendelig stor avstand / ikke oppdaget ennå), og settes først når en node oppdages. Hvis dette aldri skjer er noden uoppnåelig. Ved utskrift vises slike felt som «-». Dette indikerer korrekt at ingen vei eksisterer.

Determinisme: BFS gir alltid samme avstander, men hvilken node som blir forgjenger (`parent`) kan variere dersom flere nabover oppdages samtidig. Køen (`queue`) er FIFO, rekkefølgen bestemmes av hvordan naboen ligger i `adj[u]`. For stabil og forutsigbar rekkefølge sorteres nabolistene før søket:

```
1 for(auto& nbrs : adj) sort(nbrs.rbegin(), nbrs.rend());
```

Listing 2: Rekkefølgen på utskriften (synkende)

Dette påvirker ikke avstandene, kun utskriftsrekkefølgen.

Skalerbarhet: `vector<vector<int>` fungerer godt for små og mellomstore grafer. For svært store grafer er Compressed Sparse Row (CSR) mer effektiv. Denne metoden har mer kompakt minnebruk (kun to arrayer: `offsets` og `to`) (se i deloppgaven 2). Den har i tillegg bedre cache-lokalitet ved gjennomgang av nabover. CSR bør brukes når grafen har millioner av noder eller kanter, slik som i `Skandinavia.txt`.

Deloppgave 2: Topologisk sortering

I denne deloppgaven valgte jeg å bruke Kahns algoritme (BFS) for topologisk sortering og CSR for store grafer (`Skandinavia.txt`). Topologisk sortering finnes kun for DAGer (Directed Acyclic Graphs)/rettet, uten sykler. Grafen representeres i et komprimert radformat (CSR) med to arrayer: `offsets` av størrelse $N+1$ og `to` av størrelse M . Under innlesing telles både inn- og ut-grader for hver node. Kahns algoritme starter med alle noder som har inn-grad 0, som legges i en `kø`. Deretter tas én node u ut av `køen` av gangen, legges til i resultatet, og alle utgående kanter $u \rightarrow v$ behandles ved å redusere

`indeg[v]` med 1 hvis den blir 0, legges v til i køen. Når alle noder er prosessert har vi en topologisk rekkefølge dersom ikke alle noder ble skrevet ut, betyr det at grafen inneholder en sykel (Ikke DAG).

```

1 vector<int> topo_kahn(const vector<int>& offsets, const
2   vector<int>& to, int N, vector<int> indeg) {
3   queue<int> q;
4   for(int u = 0; u < N; u++) if(indeg[u] == 0) q.push(u);
5   vector<int> order;
6   order.reserve(N);
7   while (!q.empty()) {
8     int u = q.front(); q.pop();
9     order.push_back(u);
10    for (int i = offsets[u]; i < offsets[u+1]; ++i) {
11      int v = to[i];
12      if (--indeg[v] == 0) q.push(v);
13    }
14    if((int)order.size() != N) return {};// cycle
15  }
16 }
```

Listing 3: Kahns algoritme

Kompleksitet

Vi representerer grafen i et komprimert radformat (CSR) med to arrayer: `offsets` (størrelse $N+1$) og `to` (størrelse M). Bredde-først søk (BFS) prosesserer hver node og hver kant høyst én gang:

Tidskompleksitet er $O(N + M)$. Dette er fordi hver node settes i kø én gang, og for hver node traverseres nøyaktig intervallet `to[offsets[u]..offsets[u+1]]`, som totalt summerer til M kanter. Plasskompleksitet er $O(N + M)$ for grafen (`offsets` og `to`), samt $O(N)$ tillegg for `dist`, `parent` og `queue`, dermed $O(N + M)$. Lesing av inngangsfilen er $O(N + M)$ (lineær i antall noder og kanter). CSR bygges ved prefikssummer ($O(N)$) og én passering over alle kanter ($O(M)$).

```

1 // Teller utgrad/inngrad for å bygge CSR
2 vector<int> outdeg(N, 0), indeg(N, 0);
3     for(int i = 0; i < M; i++) {
4         int u, v;
5         if(!(file >> u >> v)) {
6             cout << "Forventet" << M << "kanter, men fant
7                 f rre.\n";
8             return 1;
9         }
10        if(u < 0 || u >= N || v < 0 || v >= N) {
11            cout << "Kant utenfor omr de p linje " << (i
12                + 2) << ".\n";
13            return 1;
14        }
15        edges.emplace_back(u, v);
16        outdeg[u]++;
17        indeg[v]++;
18    }
19
20    // Bygger offsets med indekser fra CSR
21    vector<int> offsets(N+1, 0);
22    for(int u = 0; u < N; u++) offsets[u+1] = offsets[u] +
23        outdeg[u];
24
25    vector<int> cursor = offsets;
26    vector<int> to(M);
27    for(const auto& e : edges) {
28        int u = e.first; int v = e.second;
29        to[cursor[u]++] = v;
30    }
31    edges.clear(); edges.shrink_to_fit();

```

Listing 4: Noen bruk av CSR i koden

Begrensninger

Filformat: Programmet forventer første linje N M , etterfulgt av M par u v . Noder må ligge i intervallet $[0, N-1]$. For urettet graf må begge retninger stå i filen.

Minne: CSR krever to arrayer av størrelsesorden N og M . For svært store grafer kan minne bli flaskehals. Vi bruker `int` for indekser; dersom N eller M kan overstige $2^{31}-1$, må 64-bits heltall (`int64_t`) vurderes.

Uoppnåelige noder: `dist[v] = -1` markerer at noden ikke nås fra start. Disse nodeiene får ingen forgjenger (`parent[v] = -1`).

Utskrift: Å skrive ut én linje per node kan bli tregt for store N . I praksis bør man begrense antall linjer som skrives (f.eks. de første 1000) eller skrive til fil.

Determinisme: BFS-avstander er entydige, men forgjenger kan variere dersom flere naboer oppdages i samme nivå. Med CSR og uendret kantrekkefølge blir resultatet deterministisk.

Retningssans: Grafen behandles som rettet. Hvis dataene logisk sett er urettede, må begge retninger være eksplisitt til stede, ellers blir avstandene feil.

Inndata-kvalitet: Programmet sjekker at alle kanter i filen er gyldige. Hver node-ID må ligge mellom 0 og $N-1$; hvis en kant peker til en node utenfor dette området, stoppes programmet med en feilmelding. Det kontrolleres også at filen faktisk inneholder alle de M kantene som er oppgitt på første linje. Dersom grafen bruker 1-baserte noder (starter fra 1 i stedet for 0), må ID-ene reduseres med 1 før grafen bygges.