Tópicos de Análise de Algoritmos

Algoritmos de clustering

Sec 4.5, 4,6 e 4.7 do KT

U: conjunto de n itens $\{p_1, \ldots, p_n\}$.

 $d(p_i, p_j)$: distância entre p_i e p_j .

- $d(p_i, p_i) = 0$ para i = 1, ..., n

U: conjunto de n itens $\{p_1, \ldots, p_n\}$.

 $d(p_i, p_j)$: distância entre p_i e p_j .

- $d(p_i, p_i) = 0 \text{ para } i = 1, \dots, n$
- $d(p_i, p_j) = d(p_j, p_i) > 0 para i \neq j$

k: número de grupos.

k-clustering: partição de U em k partes.

U: conjunto de n itens $\{p_1, \ldots, p_n\}$.

 $d(p_i, p_j)$: distância entre p_i e p_j .

- $d(p_i, p_i) = 0 \text{ para } i = 1, \dots, n$
- $d(p_i, p_j) = d(p_j, p_i) > 0 para i \neq j$

k: número de grupos.

k-clustering: partição de U em k partes.

espaçamento de *k*-clustering: distância mínima entre dois elementos de partes distintas do clustering.

U: conjunto de n itens $\{p_1, \ldots, p_n\}$.

 $d(p_i, p_j)$: distância entre p_i e p_j .

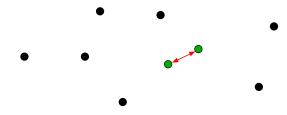
- $d(p_i, p_i) = 0 \text{ para } i = 1, \dots, n$
- $d(p_i, p_j) = d(p_j, p_i) > 0 para i \neq j$

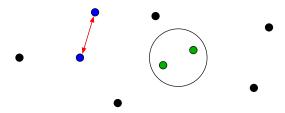
k: número de grupos.

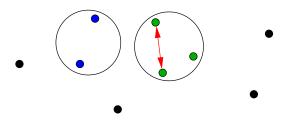
k-clustering: partição de U em k partes.

espaçamento de k-clustering: distância mínima entre dois elementos de partes distintas do clustering.









Lembram do algoritmo de Kruskal?

Problema: Dado um grafo G = (V, E) conexo, e um custo $c_e > 0$ para cada aresta e em E, encontrar uma árvore geradora mínima.

Árvore geradora mínima: árvore geradora de custo mínimo.

Lembram do algoritmo de Kruskal?

```
Problema: Dado um grafo G = (V, E) conexo, e um custo c_e > 0 para cada aresta e em E, encontrar uma árvore geradora mínima.
```

Árvore geradora mínima: árvore geradora de custo mínimo.

```
\begin{array}{lll} \mathsf{MST\text{-}KRUSKAL} \; (V, {\color{red} E}, c) \\ 1 & {\color{red} e_1, \dots, e_m} \leftarrow \mathsf{ORDENE}({\color{blue} E}, c) \\ 2 & {\color{gray} T} \leftarrow \emptyset & {\color{gray} i} \leftarrow 0 \\ 3 & \mathsf{enquanto} \; {\color{gray} G'} = (V, T) \; \mathsf{n\~ao} \; \mathsf{\'e} \; \mathsf{conexo} \; \mathsf{faça} \\ 4 & {\color{gray} i} \leftarrow {\color{gray} i} + 1 \\ 5 & \mathsf{se} \; {\color{gray} T} \cup \{{\color{gray} e_i}\} \; \mathsf{n\~ao} \; \mathsf{tem} \; \mathsf{circuito} \\ 6 & {\color{gray} ent\~ao} \; {\color{gray} T} \leftarrow {\color{gray} T} \cup \{{\color{gray} e_i}\} \\ 7 & \mathsf{devolva} \; {\color{gray} T} \end{array}
```

Como adaptar MST-KRUSKAL para resolver o k-clustering?

Como adaptar MST-KRUSKAL para resolver o k-clustering?

Acrescentamos arestas até termos k componentes.

Como adaptar MST-KRUSKAL para resolver o k-clustering?

Acrescentamos arestas até termos k componentes.

```
CLUSTERING-KRUSKAL (V, E, c, k)

1 e_1, \ldots, e_m \leftarrow \mathsf{ORDENE}(E, c) \triangleright pelo valor de c_e

2 F \leftarrow \emptyset i \leftarrow 0

3 enquanto G' = (V, F) tem > k componentes faça

4 i \leftarrow i + 1

5 se F \cup \{e_i\} não tem circuito

6 então F \leftarrow F \cup \{e_i\}

7 devolva F
```

Como adaptar MST-KRUSKAL para resolver o k-clustering?

Acrescentamos arestas até termos k componentes.

```
CLUSTERING-KRUSKAL (V, E, c, k)

1 e_1, \dots, e_m \leftarrow \mathsf{ORDENE}(E, c) \triangleright pelo valor de c_e

2 F \leftarrow \emptyset i \leftarrow 0

3 enquanto G' = (V, F) tem > k componentes faça

4 i \leftarrow i + 1

5 \mathsf{se}\ F \cup \{e_i\} não tem circuito

6 então F \leftarrow F \cup \{e_i\}

7 devolva F
```

Este algoritmo resolve o problema?

```
CLUSTERING-KRUSKAL (V, E, c, k)

1 e_1, \ldots, e_m \leftarrow \mathsf{ORDENE}(E, c) \triangleright pelo valor de c_e

2 F \leftarrow \emptyset i \leftarrow 0

3 enquanto G' = (V, F) tem \geq k componentes faça

4 i \leftarrow i + 1

5 \mathsf{se}\ F \cup \{e_i\} não tem circuito

6 então F \leftarrow F \cup \{e_i\}

7 devolva as componentes de F
```

Qual é o espaçamento do clustering devolvido?

```
CLUSTERING-KRUSKAL (V, E, c, k)

1 e_1, \dots, e_m \leftarrow \mathsf{ORDENE}(E, c) \triangleright pelo valor de c_e

2 F \leftarrow \emptyset i \leftarrow 0

3 enquanto G' = (V, F) tem \geq k componentes faça

4 i \leftarrow i + 1

5 se F \cup \{e_i\} não tem circuito

6 então F \leftarrow F \cup \{e_i\}

7 devolva as componentes de F
```

Qual é o espaçamento do clustering devolvido?

É o custo da próxima aresta que seria incluída em F.

```
CLUSTERING-KRUSKAL (V, E, c, k)

1 e_1, \ldots, e_m \leftarrow \mathsf{ORDENE}(E, c) \triangleright pelo valor de c_e

2 F \leftarrow \emptyset i \leftarrow 0

3 enquanto G' = (V, F) tem \geq k componentes faça

4 i \leftarrow i + 1

5 se F \cup \{e_i\} não tem circuito

6 então F \leftarrow F \cup \{e_i\}

7 devolva as componentes de F
```

Qual é o espaçamento do clustering devolvido?

É o custo da próxima aresta que seria incluída em F.

A próxima aresta que não forma circuito com *F*.



Seja e a próxima aresta que seria incluída em F.

Seja *e* a próxima aresta que seria incluída em *F*.

```
Seja \mathcal{C} = \{C_1, \dots, C_k\} o k-clustering devolvido, e \mathcal{C}^* = \{C_1^*, \dots, C_k^*\} um k-clustering ótimo.
```

Seja *e* a próxima aresta que seria incluída em *F*.

Seja
$$C = \{C_1, \dots, C_k\}$$
 o k -clustering devolvido, e $C^* = \{C_1^*, \dots, C_k^*\}$ um k -clustering ótimo.

O espaçamento de \mathcal{C} é c_e .

Basta mostrar que o espaçamento de \mathcal{C}^* é no máximo c_e .

Seja e a próxima aresta que seria incluída em F.

Seja
$$\mathcal{C} = \{C_1, \dots, C_k\}$$
 o k -clustering devolvido, e $\mathcal{C}^* = \{C_1^*, \dots, C_k^*\}$ um k -clustering ótimo.

O espaçamento de \mathcal{C} é c_e .

Basta mostrar que o espaçamento de \mathcal{C}^* é no máximo c_e .

Ou $\mathcal{C} = \mathcal{C}^*$, ou algum C_i não está contido em nenhum C_j^* .

Seja e a próxima aresta que seria incluída em F.

Seja
$$\mathcal{C} = \{C_1, \dots, C_k\}$$
 o k -clustering devolvido, e $\mathcal{C}^* = \{C_1^*, \dots, C_k^*\}$ um k -clustering ótimo.

O espaçamento de \mathcal{C} é c_e .

Basta mostrar que o espaçamento de \mathcal{C}^* é no máximo c_e .

Ou $\mathcal{C}=\mathcal{C}^*$, ou algum C_i não está contido em nenhum C_j^* . Seja j tal que $C_i\cap C_j^*\neq\emptyset$.

Seja e a próxima aresta que seria incluída em F.

Seja
$$C = \{C_1, \dots, C_k\}$$
 o k -clustering devolvido, e $C^* = \{C_1^*, \dots, C_k^*\}$ um k -clustering ótimo.

O espaçamento de \mathcal{C} é c_e .

Basta mostrar que o espaçamento de \mathcal{C}^* é no máximo c_e .

Ou $C = C^*$, ou algum C_i não está contido em nenhum C_i^* .

Seja j tal que $C_i \cap C_i^* \neq \emptyset$.

Seja f aresta de F em C_i com um único extremo em C_j^* .



Seja e a próxima aresta que seria incluída em F.

Seja
$$C = \{C_1, \dots, C_k\}$$
 o k -clustering devolvido, e $C^* = \{C_1^*, \dots, C_k^*\}$ um k -clustering ótimo.

O espaçamento de \mathcal{C} é c_e .

Basta mostrar que o espaçamento de \mathcal{C}^* é no máximo c_e .

Ou $\mathcal{C} = \mathcal{C}^*$, ou algum C_i não está contido em nenhum C_j^* .

Seja j tal que $C_i \cap C_i^* \neq \emptyset$.

Seja f aresta de F em C_i com um único extremo em C_j^* .

Claro que $c_f < c_e$ e espaçamento de C^* é no máximo c_f .



Como se implementa este algoritmo?

Como se implementa este algoritmo?

Como se implementa o algoritmo de Kruskal?

Como se implementa este algoritmo?

Como se implementa o algoritmo de Kruskal?

Estrutura de dados para conjuntos disjuntos...

Como se implementa este algoritmo?

Como se implementa o algoritmo de Kruskal?

Estrutura de dados para conjuntos disjuntos...

makeset(x)

- 1. $pai[x] \leftarrow x$
- 2. $rank[x] \leftarrow 0$

> heurística dos tamanhos

- 1. se $x \neq pai[x]$
- 2. **então** $pai[x] \leftarrow findset(pai[x]) > compressão de caminhos$

3. devolva pai[x]

```
findset(x)
makeset(x)
                                 1. se x \neq pai[x]
1. pai[x] \leftarrow x
2. rank[x] \leftarrow 0
                                 2. então pai[x] \leftarrow findset(pai[x])
                                 3. devolva pai[x]
                                 link(x, y)
union(x, y)
                                 1. se rank[x] > rank[y]
1. x' \leftarrow \mathsf{findset}(x)
                                 2. então pai[y] \leftarrow x
2. y' \leftarrow findset(y)
                                 3. senão pai[x] \leftarrow y
3. se x' \neq y'
                                 4.
                                                 se rank[x] = rank[y]
4. então link(x', y')
                                                      então rank[y] \leftarrow rank[y] + 1
```

```
findset(x)
makeset(x)
                                 1. se x \neq pai[x]
1. pai[x] \leftarrow x
2. rank[x] \leftarrow 0
                                 2. então pai[x] \leftarrow findset(pai[x])
                                 3. devolva pai[x]
                                 link(x, y)
union(x, y)
                                 1. se rank[x] > rank[y]
1. x' \leftarrow \mathsf{findset}(x)
                                 2. então pai[y] \leftarrow x
2. y' \leftarrow findset(y)
                                 3. senão pai[x] \leftarrow y
3. se x' \neq y'
                                 4.
                                                 se rank[x] = rank[y]
4. então link(x', y')
                                                      então rank[y] \leftarrow rank[y] + 1
```

Altura da árvore enraizada em x: comprimento de um caminho mais longo de x até um descendente de x na árvore.

Note que a altura da árvore enraizada em x é no máximo rank[x].

```
findset(x)
makeset(x)
                                  1. se x \neq pai[x]
1. pai[x] \leftarrow x
                                          então pai[x] \leftarrow findset(pai[x])
2. rank[x] \leftarrow 0
                                  3. devolva pai[x]
                                  link(x, y)
union(x, y)
                                  1. se rank[x] > rank[y]
1. x' \leftarrow \mathsf{findset}(x)
                                  2. então pai[y] \leftarrow x
2. y' \leftarrow findset(y)
                                  3. senão pai[x] \leftarrow y
3. se x' \neq y'
                                 4.
                                                  se rank[x] = rank[y]
4. então link(x', y')
                                                       então rank[y] \leftarrow rank[y] + 1
```

Altura da árvore enraizada em x: comprimento de um caminho mais longo de x até um descendente de x na árvore.

Note que a altura da árvore enraizada em x é no máximo rank[x].

Sempre que pendura uma subárvore em y, se for de mesma altura, aumenta o valor de rank[y].



```
findset(x)
makeset(x)
                                 1. se x \neq pai[x]
1. pai[x] \leftarrow x
                                 2. então pai[x] \leftarrow findset(pai[x])
2. rank[x] \leftarrow 0
                                 3. devolva pai[x]
                                 link(x, y)
union(x, y)
                                 1. se rank[x] > rank[y]
1. x' \leftarrow findset(x)
                                        então pai[y] \leftarrow x
2. y' \leftarrow findset(y)
                                 3. senão pai[x] \leftarrow y
3. se x' \neq y'
                                 4.
                                                 se rank[x] = rank[y]
4. então link(x', y')
                                 5.
                                                      então rank[y] \leftarrow rank[y] + 1
```

Altura da árvore enraizada em x: comprimento de um caminho mais longo de x até um descendente de x na árvore.

Note que a altura da árvore enraizada em x é no máximo rank[x].

Exercício: Mostre que se x é uma raiz e rank[x] = r, então a árvore enraizada em x tem pelo menos 2^r elementos.



Union-Find: análise de pior caso

```
makeset(x)
                                findset(x)
                                1. se x \neq pai[x]
1. pai[x] \leftarrow x
                                2. então pai[x] \leftarrow findset(pai[x])
2. rank[x] \leftarrow 0
                                3. devolva pai[x]
                                link(x, y)
union(x, y)
                                1. se rank[x] > rank[y]
1. x' \leftarrow \mathsf{findset}(x)
                                2. então pai[y] \leftarrow x
                       3. senão pai[x] \leftarrow y
2. v' \leftarrow findset(v)
3. se x' \neq y'
                                4. se rank[x] = rank[y]
4. então link(x', y')
                                                     então rank[y] \leftarrow rank[y] + 1
```

Exercício: Mostre que se x é uma raiz e rank[x] = r, então a árvore enraizada em x tem pelo menos 2^r elementos.

Conclua que, se n é o número de elementos nas florestas, então cada operação custa no pior caso $O(\lg n)$.



Union-Find: análise amortizada

```
findset(x)
makeset(x)
                                 1. se x \neq pai[x]
1. pai[x] \leftarrow x
2. rank[x] \leftarrow 0
                                 2. então pai[x] \leftarrow findset(pai[x])
                                 3. devolva pai[x]
                                 link(x, y)
union(x, y)
                                 1. se rank[x] > rank[y]
1. x' \leftarrow \mathsf{findset}(x)
                                 2. então pai[y] \leftarrow x
2. y' \leftarrow findset(y)
                                 3. senão pai[x] \leftarrow y
3. se x' \neq y'
                                 4.
                                                 se rank[x] = rank[y]
4. então link(x', y')
                                                      então rank[y] \leftarrow rank[y] + 1
```

Próximas aulas: análise amortizada.

Union-Find: análise amortizada

```
findset(x)
makeset(x)
                                 1. se x \neq pai[x]
1. pai[x] \leftarrow x
                                 2. então pai[x] \leftarrow findset(pai[x])
2. rank[x] \leftarrow 0
                                 3. devolva pai[x]
                                 link(x, y)
union(x, y)
                                 1. se rank[x] > rank[y]
1. x' \leftarrow \mathsf{findset}(x)
                                 2. então pai[y] \leftarrow x
2. y' \leftarrow findset(y)
                                 3. senão pai[x] \leftarrow y
3. se x' \neq y'
                                 4.
                                                 se rank[x] = rank[y]
4. então link(x', y')
                                                      então rank[y] \leftarrow rank[y] + 1
```

Próximas aulas: análise amortizada.

Vamos mostrar que cada operação acima tem custo amortizado de $O(\lg^* n)$, onde $\lg^* n$ é o número de vezes que temos que aplicar o \lg até atingir um número menor ou igual a 1.

