

Algorithms challenge: rank-sim

# Order items by similarity

In this tick, your aim is to find a good order for a set of items, given similarity scores between them. You are given a list of pairs of items and their similarity scores (this list doesn't include all pairs). Here is an example:

- [ticksim\\_train.csv](#)

We saw an illustration in the [video for section 6.6](#). We were given a list of students, and also the similarity scores between their submitted code for an Algorithms tick. We used Kruskal's algorithm to find an ordering for the students, such that two students with a high similarity score appeared close to each other in the order.

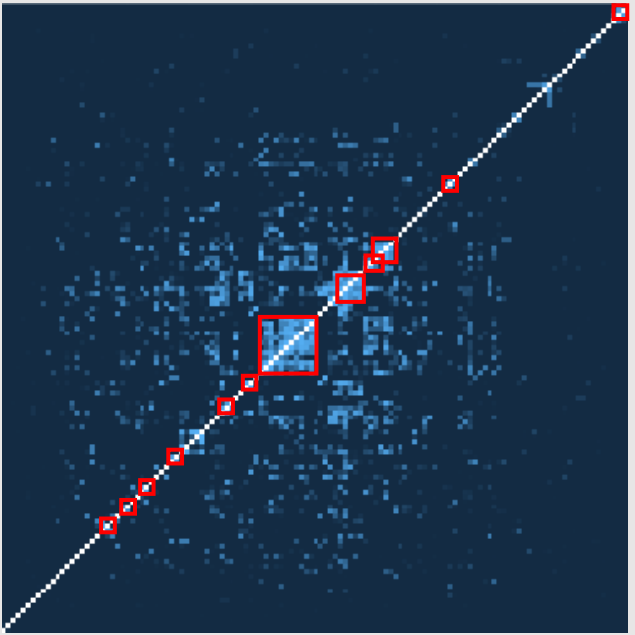
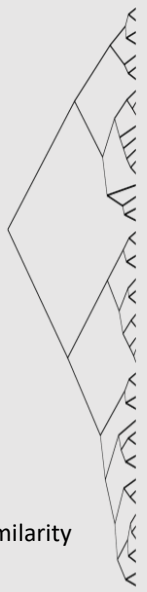
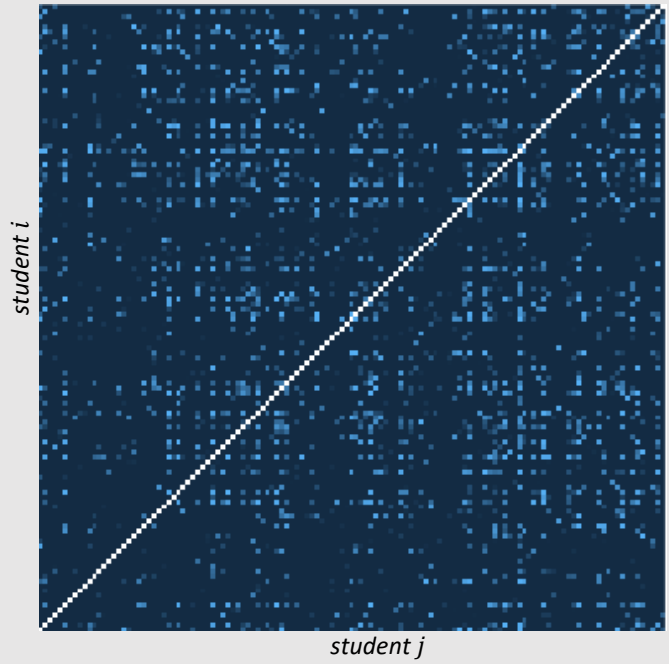
Your aim is to produce a good ordering of items. To be precise, let  $s_{uv} \in (0, 1)$  be the similarity score between items  $u$  and  $v$ . Your score will be

$$\text{score} = 100 \times \frac{x - m}{-m} \quad \text{where} \quad x = \frac{1}{MN} \sum_{\text{pairs } (u,v)} |z_u - z_v| \log(1 - s_{uv}).$$

Here  $z_u$  is the index of item  $u$  in your ordering,  $M$  is the number of pairs, and  $N$  is the number of items. The normalization is so that the score is always  $\leq 100$  (since  $x \leq 0$ ); and the constant  $m$  is the expected score from a random ordering,

$$m = \frac{1}{3M} \sum_{\text{pairs } (u,v)} \log(1 - s_{uv}),$$

Similarity matrix of max-flow code



- high similarity
- low similarity

	<b>score on training data (2021 tick1)</b>	<b>score on holdout data (2023 max-flow)</b>
Matej Urban (Trinity)	76.39	58.20
Kevin Xie (Downing)	76.46	57.39
Reuben Carolan (Churchill)	76.85	57.36
Mario Pariona Molocho (Trinity)	71.76	57.25
Wei Chuen Sin (Hughes Hall)		57.13
Leo Takashige (Trinity)	75.50	53.55
Paul DSouza (Robinson)	71.12	48.75
Ugo Obudulu (Churchill)	71.88	47.20
Katy Thackray (New Hall)	67.50	38.26
Milos Puric (Trinity)	65.41	33.97
Leonard Ong (Hughes Hall)	65.41	33.97
Elizabeth Ho (Trinity)	64.98	32.52
Dhruv Pattem (Trinity)	31.50	22.90
Max Bowman (Robinson)	0.51	1.06
George Ogden (Trinity)	76.77	0.00

# The Greatest Algorithm Anybody has ever Written

Reuben Carolan

# Optimization Problem

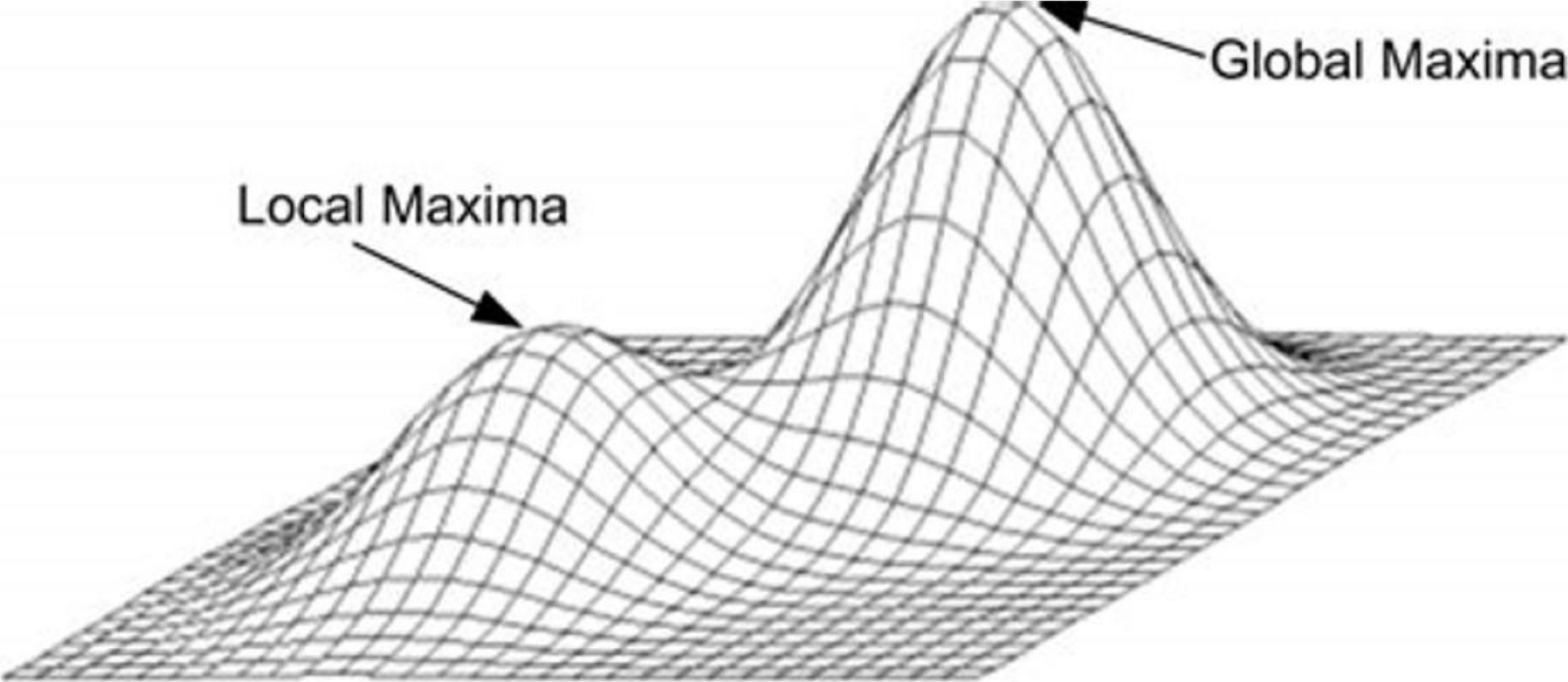
- We have function.
- We want inputs that give function biggest possible value.
  
- In this case find an order of the list to maximise:

$$\sum_{\text{pairs } (u,v)} |z_u - z_v| \log(1 - s_{uv}).$$

# Solution Space

- Solutions have neighbours. In this case solutions are neighbours if they are the same but with two items swapped places.
- We could naively randomly pick a neighbour and move to that solution if it scores better.

# Solution Space



# Simulated Annealing

- We define a Temperature value that decreases over time.
- And for each potential solution an accept probability:

```
accept_p = np.exp(( newscore - oldscore ) / curr_temperature)
```

- We Accept a new solution is better if:  
(New solution is better OR a random number from  $[0,1] < \text{accept\_p}$ )



# Why Simulated Annealing

- Idk  $\_ \_ (\_ \_) \_ / \_$
- Its good when the solution space is huge but has a small diameter. (In this case we are only ever 152 swaps from the perfect solution).
- Its good at getting a 'pretty good' solution in a short amount of time.
- It works well on the Travelling Salesman problem which is another problem of ordering a list.



# Improvements and Other Algorithms

- Restarting
- Better Candidate Generation
- Ant Colony Optimisation

SECTION X  
Optimization  
algorithms

*Non-examinable.*

# Max-flow

Given a directed graph  $g = (V, E)$  in which each edge  $u \rightarrow v$  has a capacity  $c_{uv}$ , and given a source  $s$  and a sink  $t$ , find a flow from  $s$  to  $t$  with maximum possible value.

Let  $f_{uv}$  = flow on edge  $u \rightarrow v$ .

We want to maximize flow value = net flow out of  $s$  =  $\sum_{u: s \rightarrow u} f_{su} - \sum_{w: w \rightarrow s} f_{ws}$

For it to be a valid flow, we need  $0 \leq f_{uv} \leq c_{uv}$  for all edges,

and flow conservation at every vertex  $v \in V \setminus \{s, t\}$ , i.e.  $\sum_{u: v \rightarrow u} f_{vu} - \sum_{w: w \rightarrow v} f_{wv} = 0 \quad \forall v \in V \setminus \{s, t\}$ .

Restate the problem:

maximize  $\sum_{e \in E} b_e f_e$  where  $b_e = \begin{cases} 1 & \text{if edge } e \text{ starts at } s \\ -1 & \text{if edge } e \text{ ends at } s \\ 0 & \text{otherwise} \end{cases}$

over  $f \in \mathbb{R}^E$ ,  $0 \leq f \leq c$  —  $0$  and  $c$  are vectors in  $\mathbb{R}^E$ , and the inequality must hold elementwise

such that  $\sum_{e \in E} A_{ve} f_e = 0 \quad \forall v \in V \setminus \{s, t\}$  where  $A_{ve} = \begin{cases} 1 & \text{if edge } e \text{ starts at } v \\ -1 & \text{if edge } e \text{ ends at } v \\ 0 & \text{otherwise} \end{cases}$

We can rewrite as a matrix equation

$$Af = 0$$

$$\begin{matrix} |V|-2 \\ \text{rows} \end{matrix} \left\{ \begin{matrix} |E| \text{ cols} \\ [A] \end{matrix} \right. [f] = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

```

1 import pandas
2 import numpy as np
3 df = pandas.read_csv('flownetwork_02.csv')
4 df

```

	u	v	capacity
0	0	1	10000
1	1	2	10000
2	2	5	10000
3	0	3	8000
4	3	2	8000
5	1	4	6000
6	4	5	6000

```

5 s,t = 0,5
6 V = set(df.u) | set(df.v)
7 b = np.where(df.u==s,1,0) - np.where(df.v==s,1,0)
8 A = np.row_stack([np.where(df.u==v,1,0) - np.where(df.v==v,1,0)
9                  for v in V if v not in {s,t}])
10
11 import scipy.optimize
12 res = scipy.optimize.linprog(
13     -b,
14     bounds=np.column_stack([np.zeros(len(df)), df.capacity]),
15     A_eq=A, b_eq=np.zeros(len(A))
16     )
17 df['flow'] = res.x

```

minimize  $-b \cdot x$   
 over  $x \in \mathbb{R}^n, 0 \leq x \leq \text{capacity}$   
 such that  $A_{\text{eq}}x = b_{\text{eq}}$

# Minimum spanning tree

Given a connected undirected graph where the weight of edge  $u - v$  is  $c_{uv}$ , find a spanning tree of minimum weight.

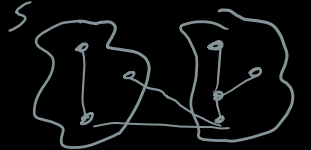
We want to find a subgraph, i.e. a subset of edges from the original graph.

Let  $x_{uv} = \begin{cases} 1 & \text{if edge } u-v \text{ is in our subgraph} \\ 0 & \text{otherwise} \end{cases}$   $x \in \{0,1\}^E$

The weight is  $\sum_{u-v} x_{uv} c_{uv}$

For  $x$  to be spanning: for every non-trivial partition of the vertices  $V = S \cup \bar{S}$ ,  $S \neq \emptyset$ ,  $S \neq V$ , there must be at least one edge between them.

Thus  $x$  spanning  $\Rightarrow \sum_{v \in S, w \in \bar{S}} x_{vw} \geq 1 \quad \forall$  non-trivial partitions  $(S, \bar{S})$ .



EXERCISE: show that if  $x$  is not spanning then  $\exists$  partition  $(S, \bar{S})$  with  $\sum_{v \in S, w \in \bar{S}} x_{vw} = 0$ .

For  $x$  to be a spanning tree: we need the spanning condition, and  $\sum_e x_e = |V| - 1$ . (EXERCISE).

Restate the problem:

minimize  $\sum_e x_e c_e$

over  $x \in \{0,1\}^E$

such that  $\sum_e x_e = |V| - 1$  and  $\sum_e A_{se} x_e \geq 1$  for all non-trivial partitions  $(S, \bar{S})$ .

This is too hard! Integer constraints are tough! Let's relax the problem, and simply require  $x \in \mathbb{R}^E$ ,  $0 \leq x \leq 1$ . Interpret  $x_e$  as a measure of "how much we want edge  $e$ ".

What other problems  
from this course can be  
written as optimization  
problems?  
as linear programs?

problems of the form

maximize  $c \cdot x$

over  $x \in \mathbb{R}^n$ ,  $\underbrace{l \leq x \leq u}_{\text{optimal}}$

such that  $Ax = b$

$A'x \leq b'$



Andrej Karpathy ✓

@karpathy



Gradient descent can write code better than you. I'm sorry.

3:56 PM - 4 Aug 2017

343 Retweets 1,161 Likes



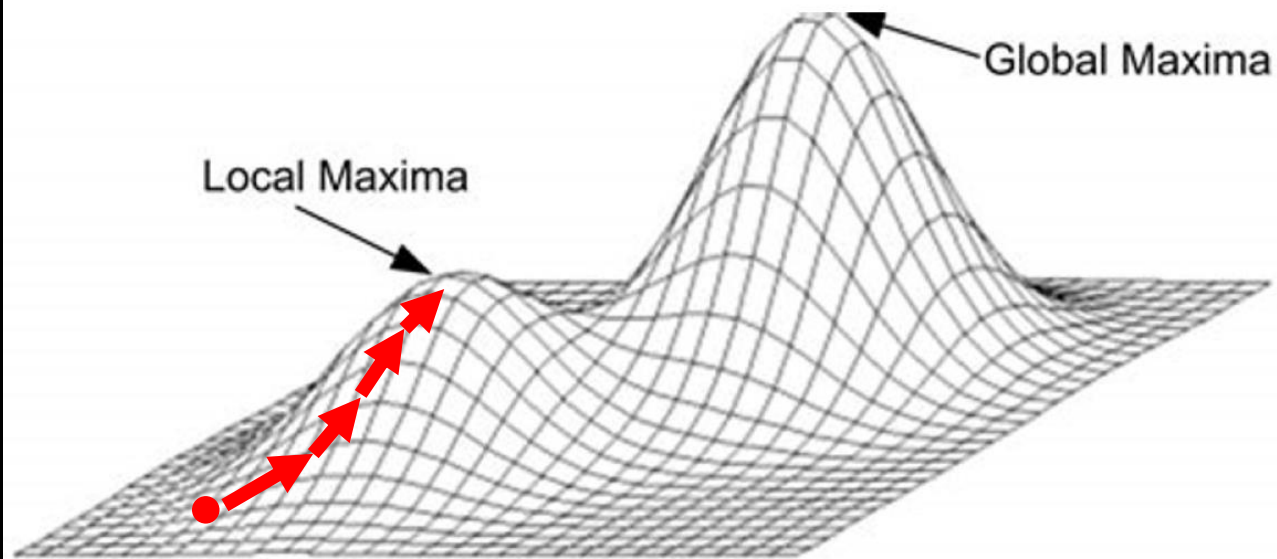
72



343



1.2K





- ❖ Can I express my task as an optimization problem?
- ❖ ... that can be solved with off-the-shelf optimizers?
- ❖ If I can't, is there an adjacent problem that's more amenable?

# Algorithms challenge: rank-sim

## Order items by similarity

In this tick, your aim is to find a good order for a set of items, given similarity scores between them. You are given a list of pairs of items and their similarity scores (this list doesn't include all pairs). Here is an example:

- [ticksim\\_train.csv](#)

We saw an illustration in the [video for section 6.6](#). We were given a list of students, and also the similarity scores between their submitted code for an Algorithms tick. We used Kruskal's algorithm to find an ordering for the students, such that two students with a high similarity score appeared close to each other in the order.

Your aim is to produce a good ordering of items. To be precise, let  $s_{uv} \in (0, 1)$  be the similarity score between items  $u$  and  $v$ . Your score will be

$$\text{score} = 100 \times \frac{x - m}{-m} \quad \text{where} \quad x = \frac{1}{MN} \sum_{\text{pairs } (u,v)} |z_u - z_v| \log(1 - s_{uv}).$$

Here  $z_u$  is the index of item  $u$  in your ordering,  $M$  is the number of pairs, and  $N$  is the number of items.

In this tick, your aim is to find a good order for a set of items, given similarity scores between them. You are given a list of pairs of items and their similarity scores (this list doesn't include all pairs). Here is an example:

- [ticksim\\_train.csv](#)

We saw an illustration in the [video for section 6.6](#). We were given a list of students, and also the similarity scores between their submitted code for an Algorithms tick. We used Kruskal's algorithm to find an ordering for the students, such that two students with a high similarity score appeared close to each other in the order.

Your aim is to produce a good ordering of items. To be precise, let  $s_{uv} \in (0, 1)$  be the similarity score between items  $u$  and  $v$ . Your score will be

$$\text{score} = 100 \times \frac{x - m}{-m} \quad \text{where} \quad x = \frac{1}{MN} \sum_{\text{pairs } (u,v)} |z_u - z_v| \log(1 - s_{uv}).$$

Here  $z_u$  is the index of item  $u$  in your ordering,  $M$  is the number of pairs, and  $N$  is the number of items.

Related (and much easier) problem:

$$\begin{aligned} &\text{maximize} && \sum_{u,v} |z_u - z_v| \log(1 - s_{uv}) \\ &\text{over} && z \in \mathbb{R}^N \\ &\text{such that} && \max(z) - \min(z) = N - 1 \end{aligned}$$

*These are called embedding problems. We want to “embed” each student into some other space*

- *either  $z_u \in \mathbb{N}$ , an integer embedding*
- *or  $z_u \in \mathbb{R}$ , a real embedding*



*The hope is that if two students are similar ( $s_{uv} \approx 1$ ) then they'll end up with similar embeddings ( $z_u \approx z_v$ ).*

# Tokenizer

The GPT family of models process text using **tokens**, which are common sequences of characters found in text. The models understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.

You can use the tool below to understand how a piece of text would be tokenized by the API, and the total count of tokens in that piece of text.

GPT-3 Codex

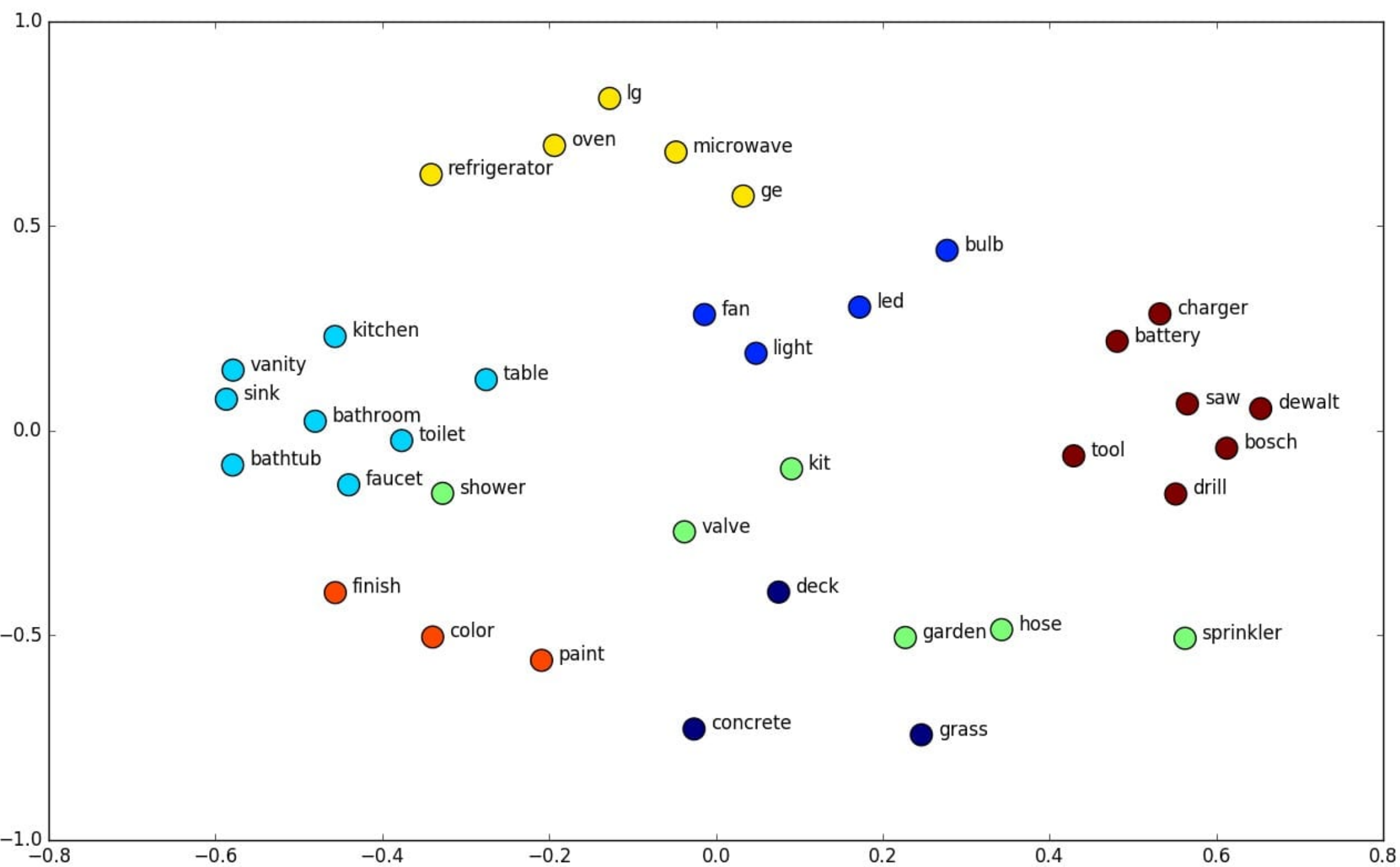
```
The renowned computer scientist Edsger Dijkstra said that "the question of whether machines can think is about as relevant as the question of whether submarines can swim."
```

Clear Show example

**Tokens**      **Characters**  
**33**          **171**

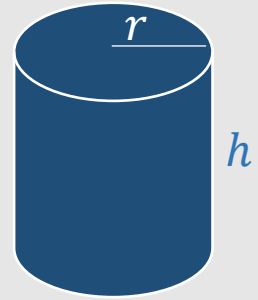
The renowned computer scientist Edsger Dijkstra said that "the question of whether machines can think is about as relevant as the question of whether submarines can swim."

TEXT      TOKEN IDS



## Maths for NST B lecture 14

A right-circular cylinder of radius  $r$  and height  $h$  has volume  $\pi r^2 h$  and surface area  $2\pi r^2 + 2\pi r h$ . Given the surface area is  $A$ , find the largest possible volume.



## Maths for NST B lecture 14

A right-circular cylinder of radius  $r$  and height  $h$  has volume  $\pi r^2 h$  and surface area  $2\pi r^2 + 2\pi r h$ . Given the surface area is  $A$ , find the largest possible volume.

$$\begin{array}{ll} \text{maximize} & \pi r^2 h \\ \text{over} & r, h \in \mathbb{R}_{\geq 0} \\ \text{such that} & 2\pi r^2 + 2\pi r h - A = 0 \end{array}$$

### METHOD

1. Write out the Lagrangian

$$\mathcal{L}(r, h; \lambda) = \pi r^2 h - \lambda(2\pi r^2 + 2\pi r h - A)$$

2. For a given  $\lambda$ , find  $r \geq 0$  and  $h \geq 0$  to maximize  $\mathcal{L}(r, h; \lambda)$
3. Choose  $\lambda$  so that these  $r$  and  $h$  satisfy the constraint



$$\text{Step 2: } \left. \begin{aligned} \frac{\partial I}{\partial r} &= 2\pi r h - \lambda (4\pi r + 2\pi h) = 0 \\ \frac{\partial I}{\partial h} &= \pi r^2 - \lambda (2\pi r) = 0 \end{aligned} \right\} \Rightarrow \begin{aligned} h &= 4\lambda \\ r &\neq 0 \text{ or } r = 2\lambda \end{aligned}$$

$$\text{Step 3: } 2\pi r^2 + 2\pi r h = A \Rightarrow 8\pi \lambda^2 + 16\pi \lambda^2 = A \Rightarrow \lambda = \sqrt{\frac{A}{24\pi}}$$

$$\text{This gives } V = \pi r^2 h = 16\pi \lambda^3 = 16\pi \left(\frac{A}{24\pi}\right)^{3/2}$$

## Maths for NST B lecture 14

A right-circular cylinder of radius  $r$  and height  $h$  has volume  $\pi r^2 h$  and surface area  $2\pi r^2 + 2\pi r h$ . Given the surface area is  $A$ , find the largest possible volume.

maximize  $\pi r^2 h$   
over  $r, h \in \mathbb{R}_{\geq 0}$   
such that  $2\pi r^2 + 2\pi r h - A = 0$

### METHOD

1. Write out the Lagrangian

$$\mathcal{L}(r, h; \lambda) = \pi r^2 h - \lambda(2\pi r^2 + 2\pi r h - A) \quad = \text{objective function}$$

2. For a given  $\lambda$ , find  $r \geq 0$  and  $h \geq 0$  to maximize  $\mathcal{L}(r, h; \lambda)$
3. Choose  $\lambda$  so that these  $r$  and  $h$  satisfy the constraint

$-\sum_i \lambda_i \text{ constraint}_i$   
one  $\lambda$  for every constraint equation

Why does this method work?

For every  $(r, h)$  such that  $\text{Area}(r, h) = A$ , and for every  $\lambda \in \mathbb{R}$ ,

$$\begin{aligned} \text{volume}(r, h) &= \pi r^2 h \quad \text{by definition of volume} \\ &= \pi r^2 h - \lambda(\text{Area}(r, h) - A) \quad \text{since } \text{Area}(r, h) = A \text{ by assumption} \\ &= \mathcal{L}(r, h; \lambda) \quad \text{by definition of } \mathcal{L} \\ &\leq \max_{r', h' \geq 0} \mathcal{L}(r', h'; \lambda) \quad \text{since we can only possibly increase } \mathcal{L} \text{ by maximizing} \\ &\quad \text{over its first two arguments} \\ &= \lambda A - 8\pi\lambda^3 \quad \text{This maximization is pretty easy, with simple calculus} \end{aligned}$$

## Maths for NST B lecture 14

A right-circular cylinder of radius  $r$  and height  $h$  has volume  $\pi r^2 h$  and surface area  $2\pi r^2 + 2\pi r h$ . Given the surface area is  $A$ , find the largest possible volume.

$$\begin{array}{ll} \text{maximize} & \pi r^2 h \\ \text{over} & r, h \in \mathbb{R}_{\geq 0} \\ \text{such that} & 2\pi r^2 + 2\pi r h - A = 0 \end{array}$$

### METHOD

1. Write out the Lagrangian

$$\mathcal{L}(r, h; \lambda) = \pi r^2 h - \lambda(2\pi r^2 + 2\pi r h - A)$$

2. For a given  $\lambda$ , find  $r \geq 0$  and  $h \geq 0$  to maximize  $\mathcal{L}(r, h; \lambda)$

3. Choose  $\lambda$  so that these  $r$  and  $h$  satisfy the constraint

*we'll get answers as functions of  $\lambda$ :  $r = r(\lambda)$ ,  $h = h(\lambda)$*

*This works in some problems, but not all. A more general method is "Find  $\lambda$  that minimizes the upper bound".*

Why does this method work?

For every  $(r, h)$  such that  $\text{Area}(r, h) = A$ , and for every  $\lambda \in \mathbb{R}$ ,

$$\begin{aligned} \text{volume}(r, h) &= \pi r^2 h \\ &= \pi r^2 h - \lambda(\text{Area}(r, h) - A) \\ &= \mathcal{L}(r, h; \lambda) \\ &\leq \max_{r', h' \geq 0} \mathcal{L}(r', h'; \lambda) \\ &= \lambda A - 8\pi \lambda^3 \end{aligned}$$

# Max-flow

Given a directed graph  $g$  in which each edge  $u \rightarrow v$  has a capacity  $c_{uv}$ , and given a source  $s$  and a sink  $t$ , find a flow from  $s$  to  $t$  with maximum possible value.

$$L(f; \lambda) = \text{value}(f) - \sum_{v \in V \setminus \{s, t\}} \lambda_v (\text{net flow}(v) - 0)$$

flow  $f$

flow is conserved

For every  $(r, h)$  such that  $\text{Area}(r, h) = A$ , and for every  $\lambda \in \mathbb{R}$ ,  $\mathbb{R}^{v-2}$

$$\text{value}(f) \text{ volume}(r, h) = \pi r^2 h$$

$$= \pi r^2 h - \lambda (\text{Area}(r, h) - A)$$

$$= L(r, h; \lambda) \quad L(f; \lambda)$$

$$\leq \max_{r', h' \geq 0} L(r', h'; \lambda) \quad \max_{0 \leq f \leq c} L(f; \lambda)$$

$$= \lambda A - 8\pi\lambda^3 = \text{capacity}(\lambda)$$

think of  $\lambda$  as a "generalized cut"

The Lagrangian is

$$\mathcal{L}(f; \lambda) = \left( \sum_{u: s \rightarrow u} f_{su} - \sum_{w: w \rightarrow s} f_{ws} \right) - \sum_{v \neq s, t} \lambda_v \left( \sum_{u: v \rightarrow u} f_{vu} - \sum_{w: w \rightarrow v} f_{wv} \right)$$

The preceding argument (called "Lagrangian weak duality") says that for any flow  $f$  and for any  $\lambda$ ,

$$\text{val}(f) \leq \max_{f': 0 \leq f' \leq C} \mathcal{L}(f'; \lambda)$$

$$\begin{aligned} \mathcal{L}(f'; \lambda) &= \sum_v \delta_v \left( \sum_{u: v \rightarrow u} f'_{vu} - \sum_{w: w \rightarrow v} f'_{wv} \right) \text{ where } \delta_v = \begin{cases} 1 & \text{if } v = s \\ 0 & \text{if } v = t \\ \lambda_v & \text{otherwise} \end{cases} \\ &= \sum_{v, u: v \rightarrow u} \delta_v f'_{vu} - \sum_{v, w: w \rightarrow v} \delta_w f'_{wv} \\ &= \sum_{a, b: a \rightarrow b} \delta_a f'_{ab} - \sum_{b, a: a \rightarrow b} \delta_b f'_{ab} \\ &= \sum_{a, b: a \rightarrow b} f'_{ab} (\delta_a - \delta_b) \end{aligned}$$

This is maximized at  $f'_{ab} = \begin{cases} c_{ab} & \text{if } \delta_a > \delta_b \\ 0 & \text{if } \delta_a < \delta_b \\ ? & \text{if } \delta_a = \delta_b \end{cases} = \begin{cases} c_{ab} & \text{if } a \in S, b \notin S \\ 0 & \text{if } a \notin S, b \in S \\ ? & \text{if } a, b \text{ on same side} \end{cases}$

Weak duality holds for any  $\lambda$ . Let's consider  $\lambda_v = \begin{cases} 1 & \text{if } v \in S \\ 0 & \text{if } v \notin S \end{cases}$  for some cut  $(S, \bar{S})$

For such a  $\lambda$ ,  $\max_{0 \leq f' \leq C} \mathcal{L}(f'; \lambda) = \sum_{a \in S, b \notin S} c_{ab}$  = capacity  $(S, \bar{S})$

There are deep connections between optimization and graph algorithms

- ❖ Starting with the Lagrangian, it is sometimes possible to derive fast algorithms, called 'primal-dual' algorithms
- ❖ Ford–Fulkerson and Kruskal are examples from this general family