

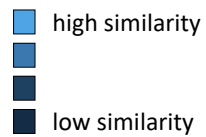
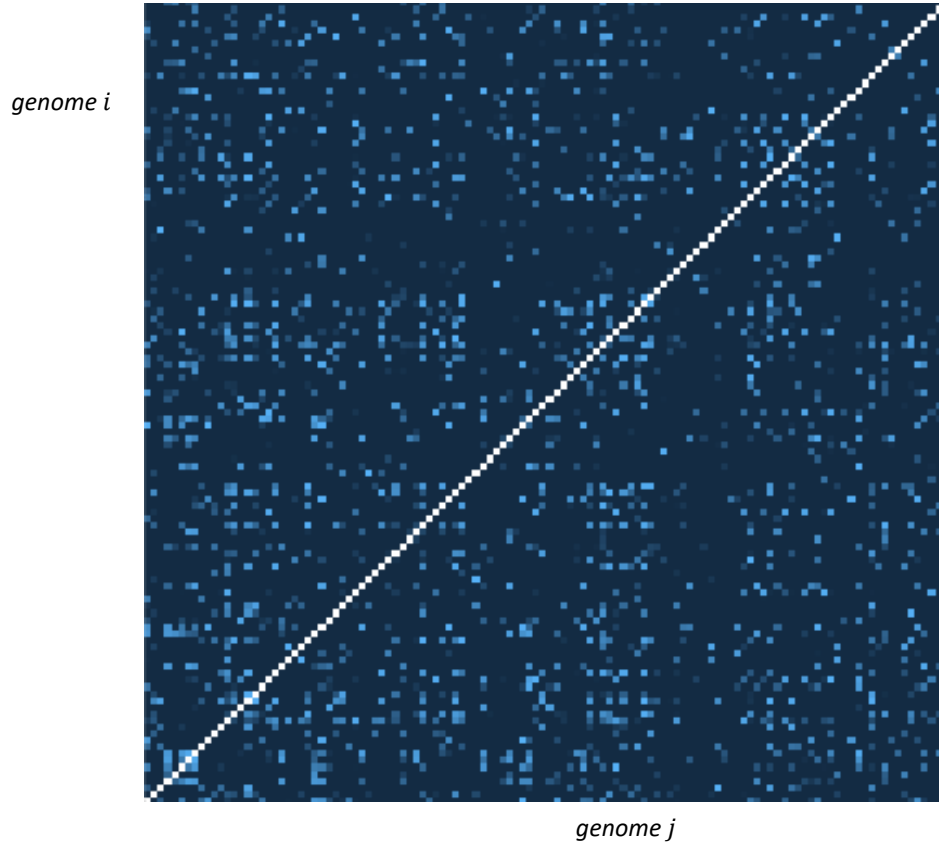
A phylogenetic tree is a tree representation of the evolutionary history of a set of gene sequences e.g. COVID variants.

SECTIONS 6.5 and 6.6

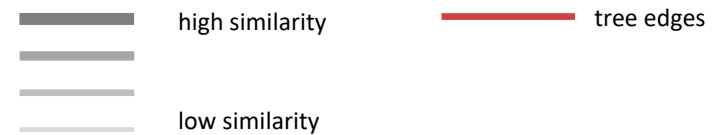
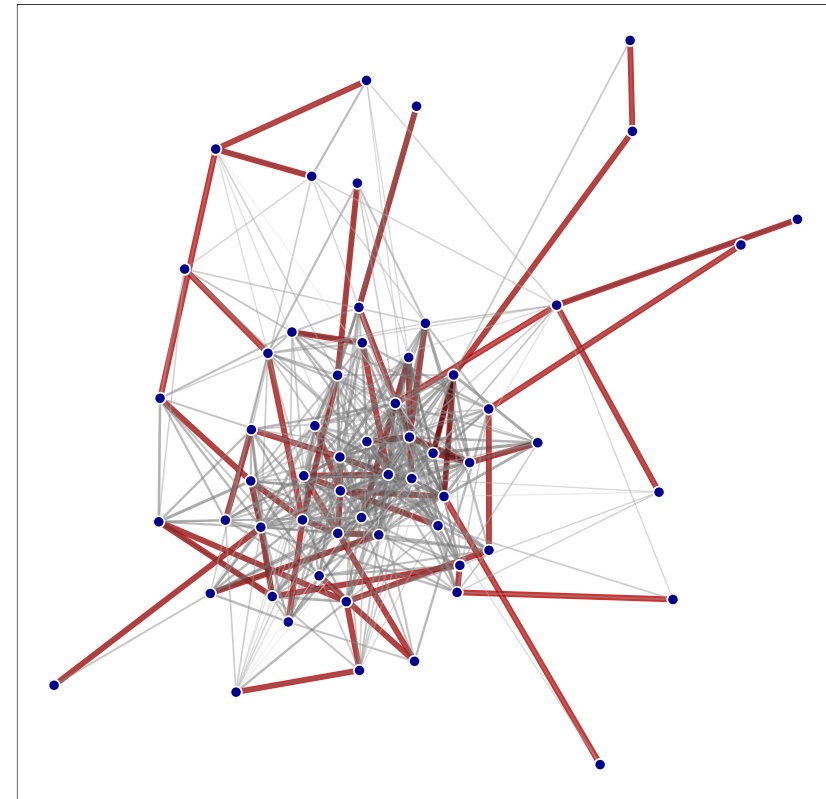
Prim's and Kruskal's algorithms

Given the similarity score between every pair of genomes, can we reconstruct a likely phylogenetic tree? In other words, can we find a high-similarity tree embedded in the similarity graph?

Similarity matrix



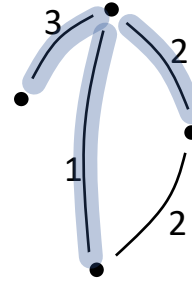
Similarity graph + subtree



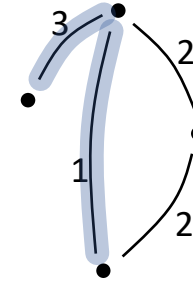
DEFINITIONS

Given a connected undirected graph g with edge weights,

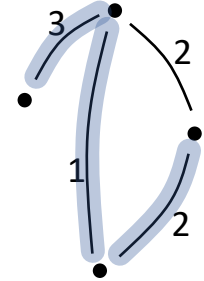
- A **spanning tree** of g is a tree that connects all of g 's vertices, using some or all of g 's edges
- The **weight** of a spanning tree is the sum of all its edge weights
- A **minimum spanning tree (MST)** is a spanning tree that has minimum weight among all spanning trees



*spanning
tree,
weight 6*



*not a
spanning tree*



*spanning
tree,
weight 6*

PROBLEM STATEMENT

Given such a graph, find a minimum spanning tree

dynamic programming

greedy algorithms

translation strategy

depth-first search

breadth-first search

Dijkstra's algorithm

Bellman-Ford algorithm

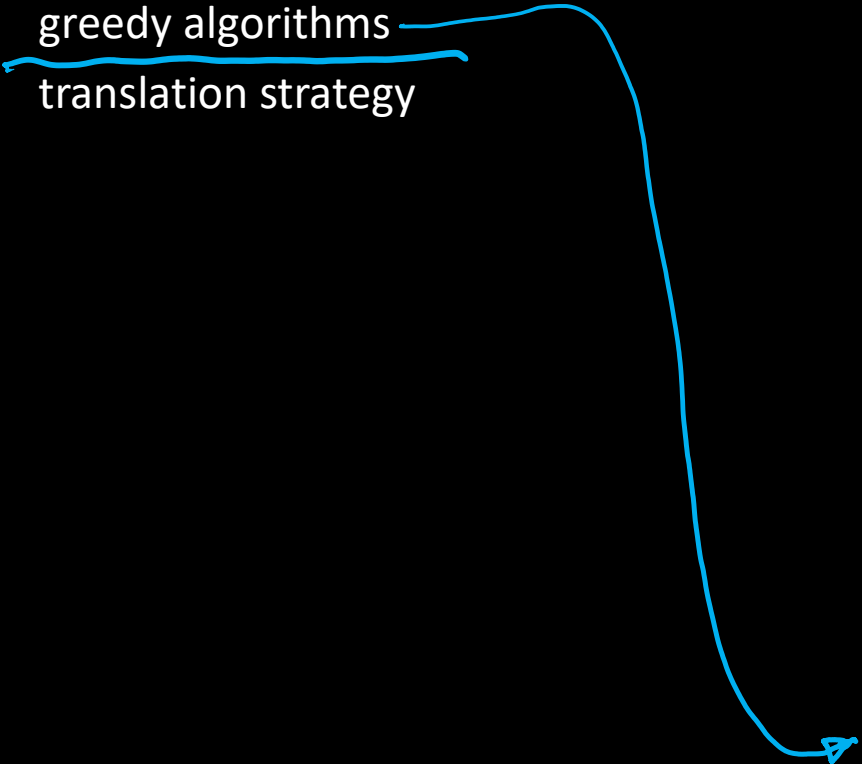
Johnson's algorithm

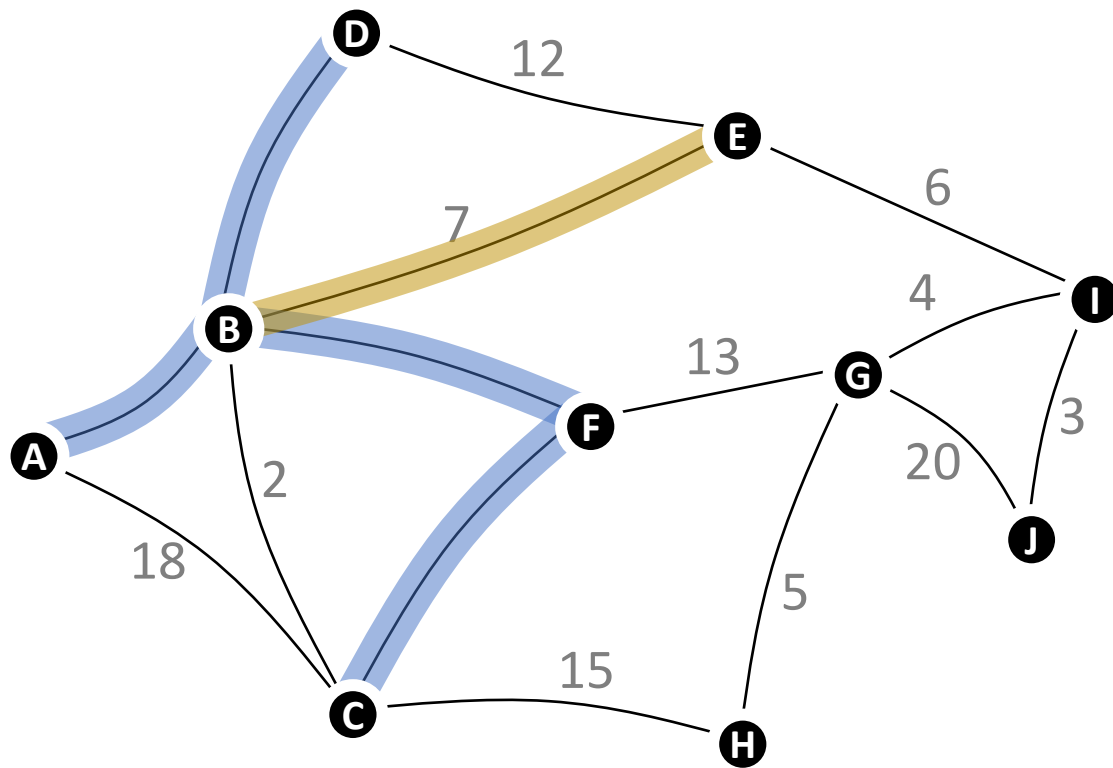
Ford-Fulkerson algorithm

matchings

topological sort

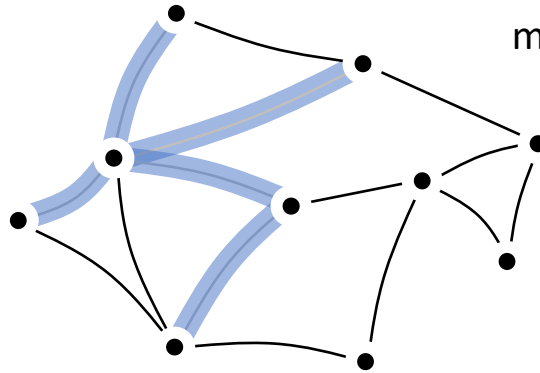
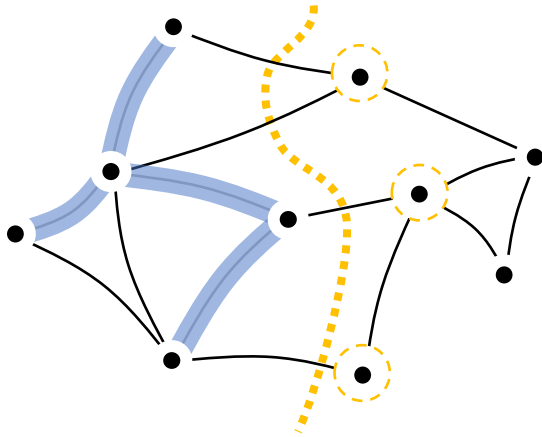
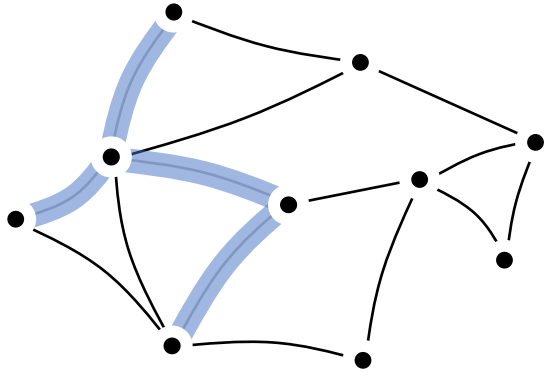
minimum spanning tree





Let's build up a **tree**, edge by edge.

SIMPLE GREEDY ALGORITHM:
Which edge would you add next, to grow the tree?



PRIM'S ALGORITHM

Choose an arbitrary start vertex as our initial tree.
Then, given a tree we've built so far,

1. look at the *frontier* of vertices we might add next, and at the *cut* between our tree and those vertices
2. pick the lowest-weight edge across this cut, and add it to the tree
3. *Assert: the tree we have so far is part of some minimum spanning tree*

Repeat until we have a spanning tree.

PROOF OF CORRECTNESS (OUTLINE)

We can prove the assertion on line 3, using the “breakpoint” proof strategy plus some fiddly maths about trees. The final output is hence a minimum spanning tree.

```

1 def prim(g, s):
2     for v in g.vertices:
3         v.distance = ∞ distance from tree to v
4         v.in_tree = False am I in the tree yet?
5     s.come_from = None
6     s.distance = 0
7     toexplore = PriorityQueue([s], sortkey = λv: v.distance)
8
9     while not toexplore.isempty():
10        v = toexplore.popmin()
11        v.in_tree = True
12        for (w, edgeweight) in v.neighbours:
13
14            if (not w.in_tree) and edgeweight < w.distance:
15                w.distance = edgeweight
16                w.come_from = v
17                if w in toexplore:
18                    toexplore.decreasekey(w)
19                else:
20                    toexplore.push(w)

```

PRIM'S ALGORITHM

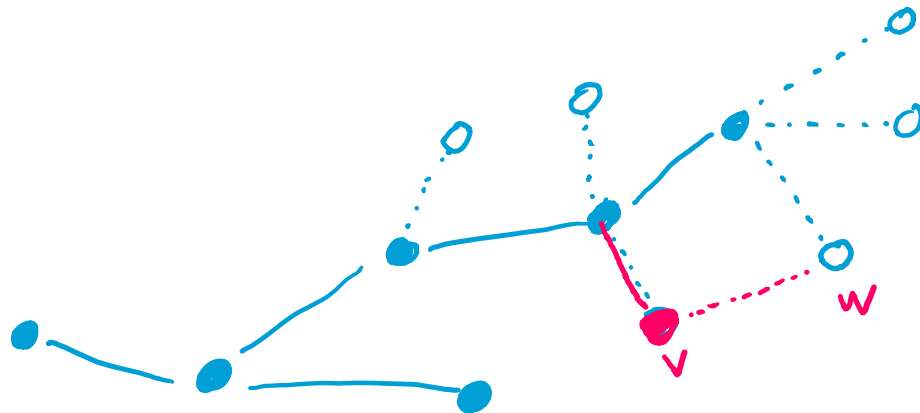
Choose an arbitrary start vertex as our initial tree.
Then, given a tree we've built so far,

1. look at the **frontier** of vertices we might add next, and at the *cut* between our tree and those vertices
2. pick the lowest-weight edge across this cut, and add it to the tree
3. *Assert: the tree we have so far is part of some minimum spanning tree*

Repeat until we have a spanning tree.

Don't recompute the frontier every iteration.

Instead, store it & update it.




```

1 def prim(g, s):
2     for v in g.vertices:
3         v.distance = ∞
4         v.in_tree = False
5     s.come_from = None
6     s.distance = 0
7     toexplore = PriorityQueue([s], sortkey = λv: v.distance)
8
9     while not toexplore.isempty():
10        v = toexplore.popmin()
11        v.in_tree = True
12        for (w, edgeweight) in v.neighbours:
13
14            if (not w.in_tree) and edgeweight < w.distance:
15                w.distance = edgeweight
16                w.come_from = v
17                if w in toexplore:
18                    toexplore.decreasekey(w)
19                else:
20                    toexplore.push(w)

```

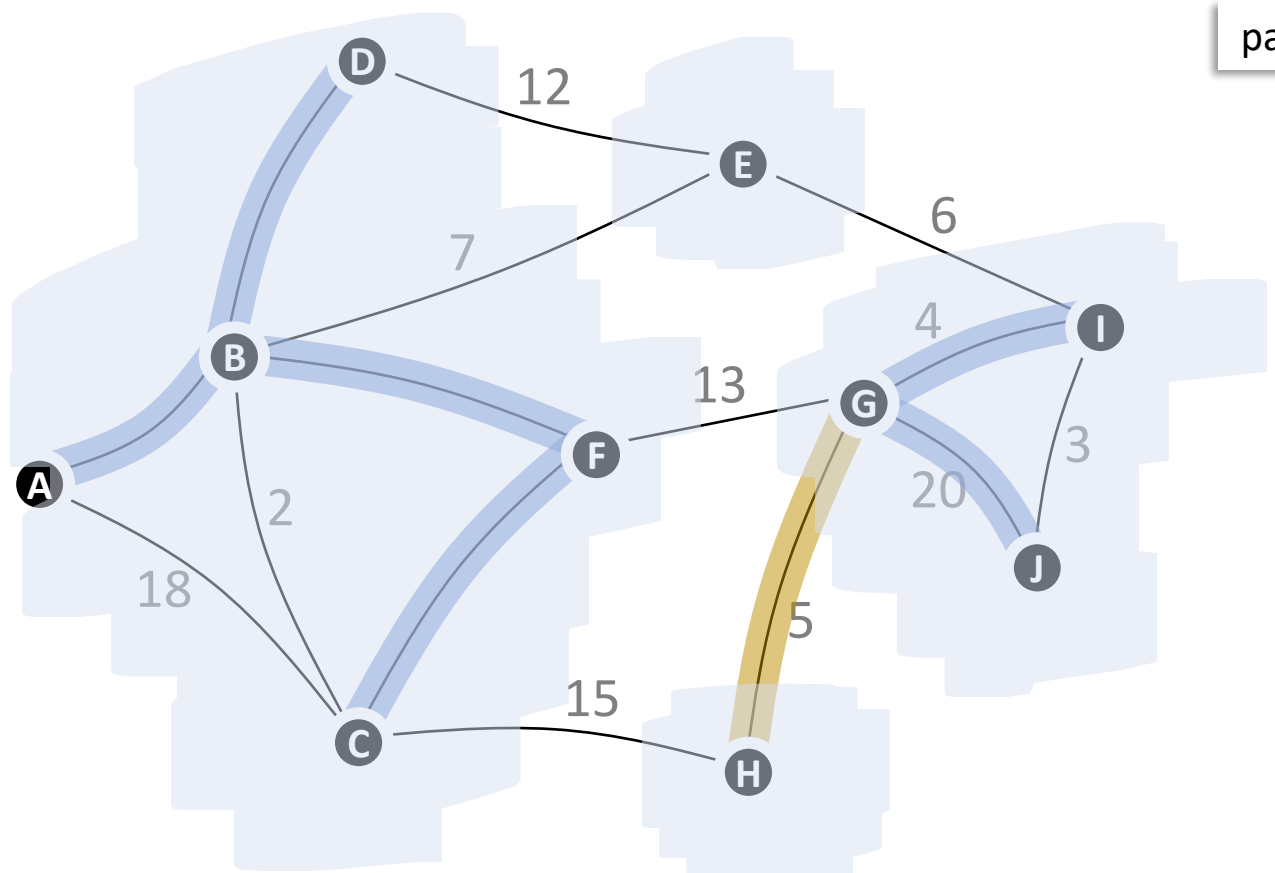
Cost $O(E + V \log V)$
because it's basically
the same as Dijkstra.

```

1 def dijkstra(g, s):
2     for v in g.vertices:
3         v.distance = ∞
4
5     s.distance = 0
6
7     toexplore = PriorityQueue([s], sortkey = λv: v.distance)
8
9     while not toexplore.is_empty():
10        v = toexplore.popmin()
11
12        for (w, edgeweight) in v.neighbours:
13            dist_w = v.distance + edgeweight
14            if dist_w < w.distance:
15                w.distance = dist + w
16
17                if w in toexplore:
18                    toexplore.decreasekey(w)
19                else:
20                    toexplore.push(w)

```

Alternatively ...
Let's build up a **forest**, edge by edge.



SIMPLE GREEDY ALGORITHM:
Which edge would you add
next, to grow the forest?

KRUSKAL'S ALGORITHM

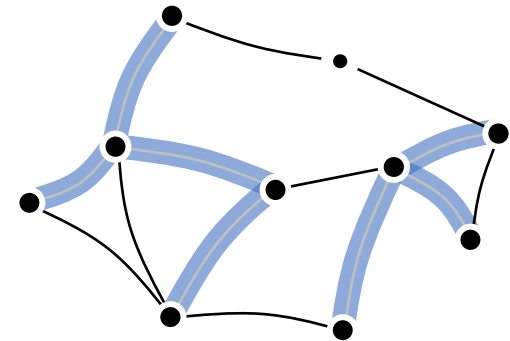
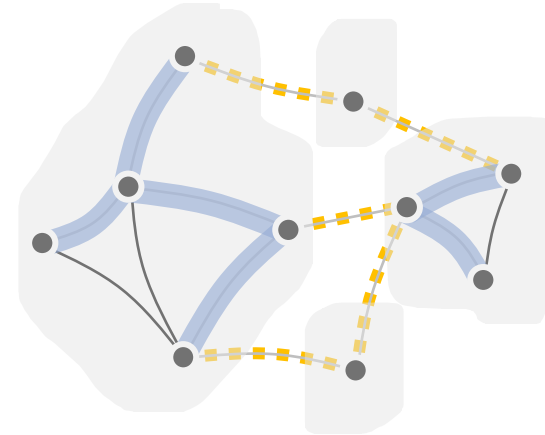
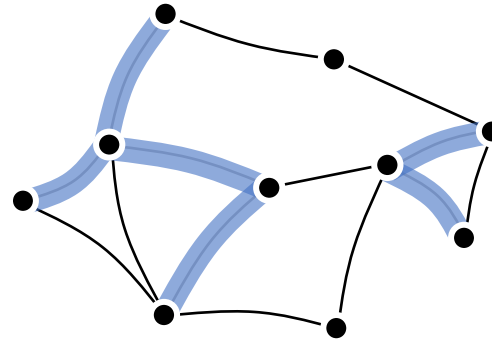
Given a **forest** we've built so far,

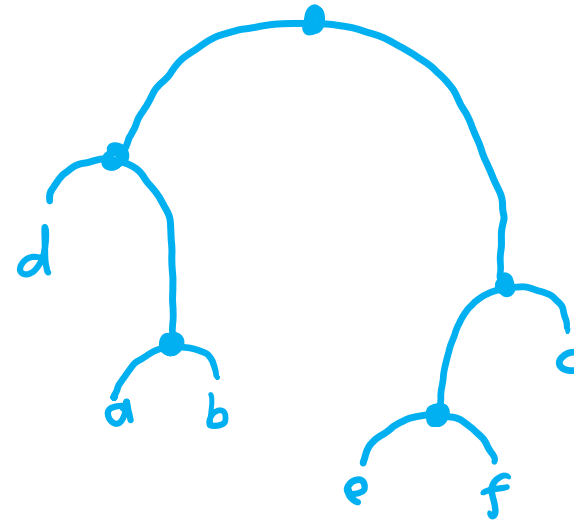
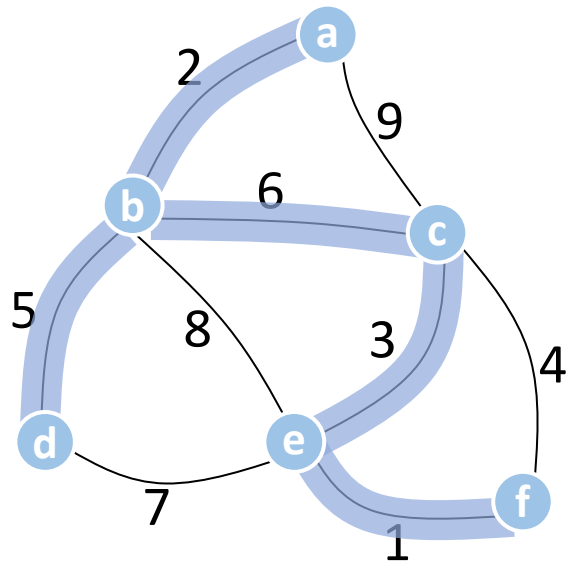
1. look at all the edges that would join two fragments of the forest
2. pick the lowest-weight one and add it to the forest, thereby joining two fragments
3. *Assert: the forest we have so far is part of some minimum spanning tree*

Repeat until we have a spanning tree.

PROOF OF CORRECTNESS (OUTLINE)

We can prove the assertion on line 3, using the “breakpoint” proof strategy plus some fiddly maths about trees. The final output is hence a minimum spanning tree.



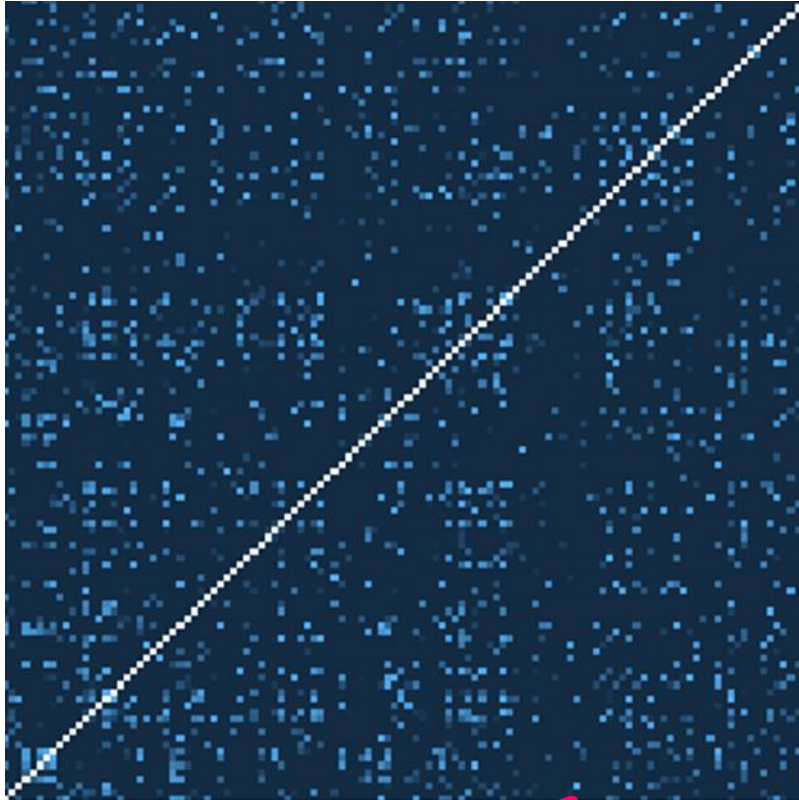


Kruskal's algorithm in effect builds a classification tree; vertices connected by low-weight edges become nearby leaves of the tree.

EXERCISE. Run through the steps of Kruskal's algorithm.

Similarity matrix of submitted coursework

~~genome i~~
student i



~~genome j~~
student j

■ high similarity
■
■
■ low similarity

Algorithms 1 2023-24 resalloc L

vle.cam.ac.uk/mod/vpl/similarity/listsimilarity.php

Dashboard Courses Find a course Categories Help About Moodle Course History

Open course index 2023-24 / resalloc

VPL **resalloc** Window Snip

Virtual programming lab Settings Test cases More

resalloc

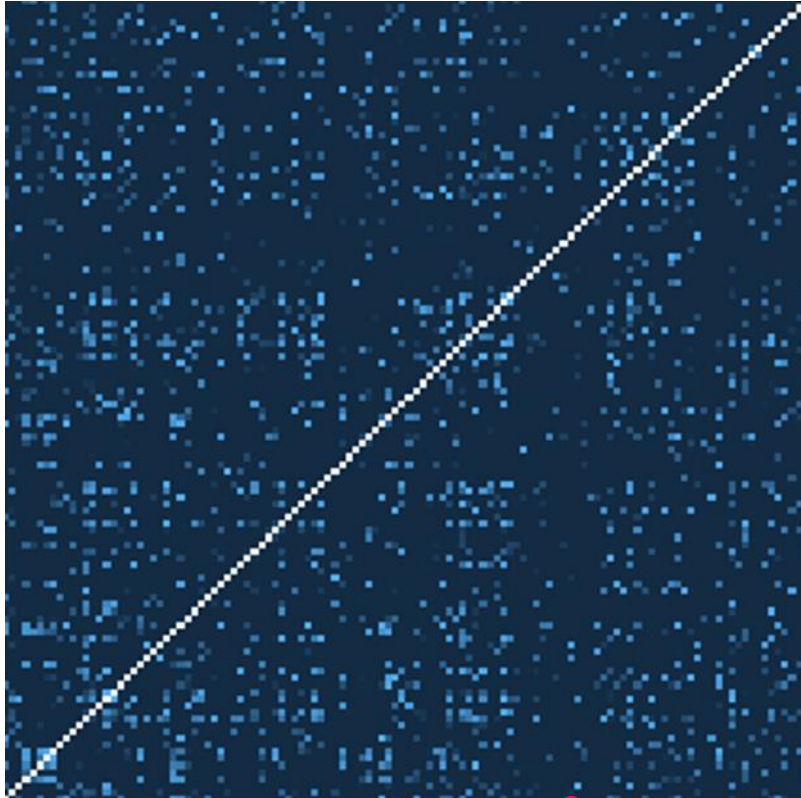
Description Submissions list **Similarity** Test activity

Similarity List of similarities found

#	First name / Last name	Similar to	Cluster #
1	resalloc.py 1.00 / 1.00 [redacted]	81 95 38*** resalloc.py 1.00 / 1.00 [redacted]	1
2	resalloc.py 1.00 / 1.00 [redacted]	74 95 24** resalloc.py 1.00 / 1.00 [redacted] (*)	2
3	resalloc.py 1.00 / 1.00 [redacted]	72 90 87*** resalloc.py 1.00 / 1.00 [redacted]	3
4	resalloc.py 1.00 / 1.00 [redacted]	80 93 60*** resalloc.py 1.00 / 1.00 [redacted]	3

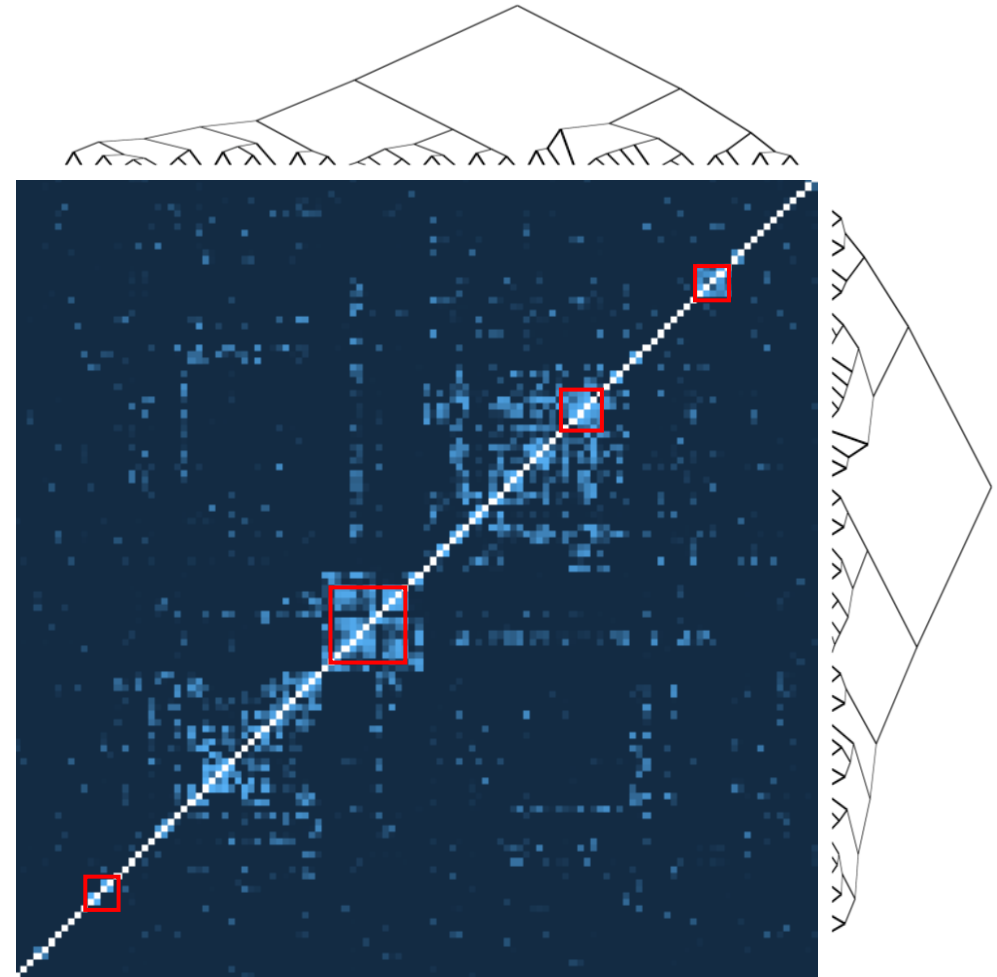
Similarity matrix of submitted coursework

~~genome i~~
student i



~~genome j~~
student j

- high similarity
- low similarity



KRUSKAL'S ALGORITHM

Given a **forest** we've built so far,

1. look at all the edges that would join two fragments of the forest
2. pick the lowest-weight one and add it to the tree, thereby joining two fragments
3. *Assert: the forest we have so far is part of some minimum spanning tree*

Repeat until we have a spanning tree.

```
1 def kruskal(g):
2     tree_edges = []
3     partition = DisjointSet()
4     for v in g.vertices:
5         partition.addsingleton(v)
6     edges = sorted(g.edges, sortkey = lambda(u,v,weight): weight)
7
8     for (u,v,edgeweight) in g.edges:
9         p = partition.getsetwith(u)
10        q = partition.getsetwith(v)
11        if p != q:
12            tree_edges.append((u,v))
13            partition.merge(p, q)
```

Don't recompute these edges every iteration.

Just pre-sort the list of all edges, then iterate through and ignore those that are within-fragment.

Total cost

$$O(V + E + E \log E)$$

$$= O(V + E \log E)$$

We're assuming a connected graph
 $\Rightarrow E \geq V - 1 \Rightarrow V \leq E + 1$

The graph can't have more than $\frac{1}{2} V(V-1)$ edges
 $\Rightarrow E \leq \frac{1}{2} V(V-1) \Rightarrow \log E \leq 2 \log V$

So total cost is $O(E \log V)$.

```

1 def kruskal(g):
2     tree_edges = []
3     partition = DisjointSet()
4     for v in g.vertices:
5         partition.addsingleton(v)
6     edges = sorted(g.edges, sortkey = lambda(u,v,weight): weight)
7
8     for (u,v,edgeweight) in g.edges:
9         p = partition.getsetwith(u)
10        q = partition.getsetwith(v)
11        if p != q:
12            tree_edges.append((u,v))
13            partition.merge(p, q)

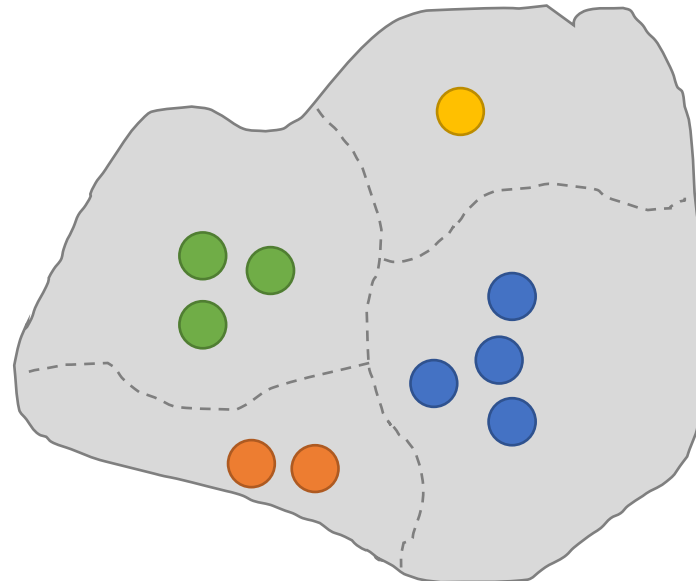
```

Annotations:

- Lines 3-5: $O(V)$
- Line 6: $O(E \log E)$
- Lines 8-13: $O(E)$ iterations
- Lines 11-13: $O(1)$ ish operations

The abstract data type **DisjointSet** stores a collection of disjoint sets, and supports

- $O(1)$ ish \blacksquare `addsingleton(v)`
- $O(1)$ ish \blacksquare `p = getsetwith(v)`
- $O(1)$ ish \blacksquare `merge(p,q)`



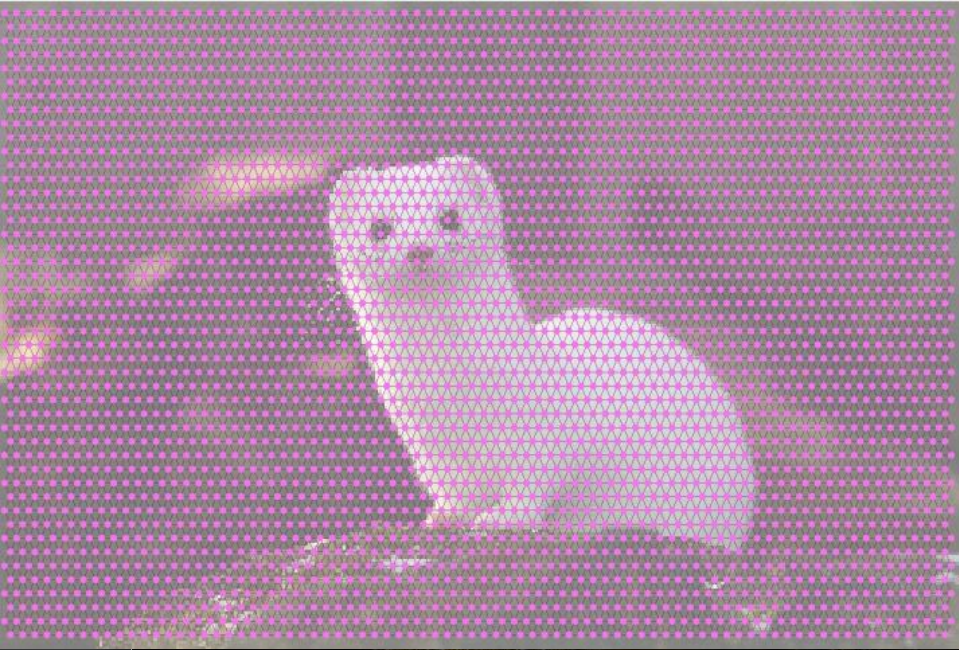


QUESTION. How might we segment this image into “handsome stoat” and “background”?

dynamic programming
greedy algorithms
translation strategy

depth-first search
breadth-first search
Dijkstra’s algorithm
Bellman-Ford algorithm
Johnson’s algorithm

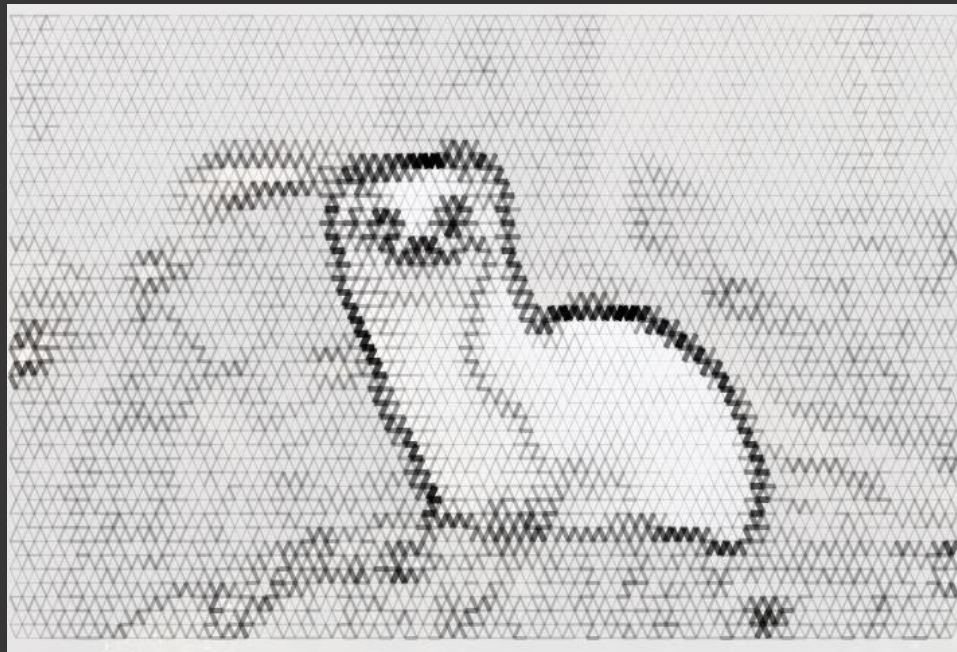
Ford-Fulkerson algorithm
matchings
topological sort
Prim, Kruskal



1. define a grid



2. measure dissimilarity along edges

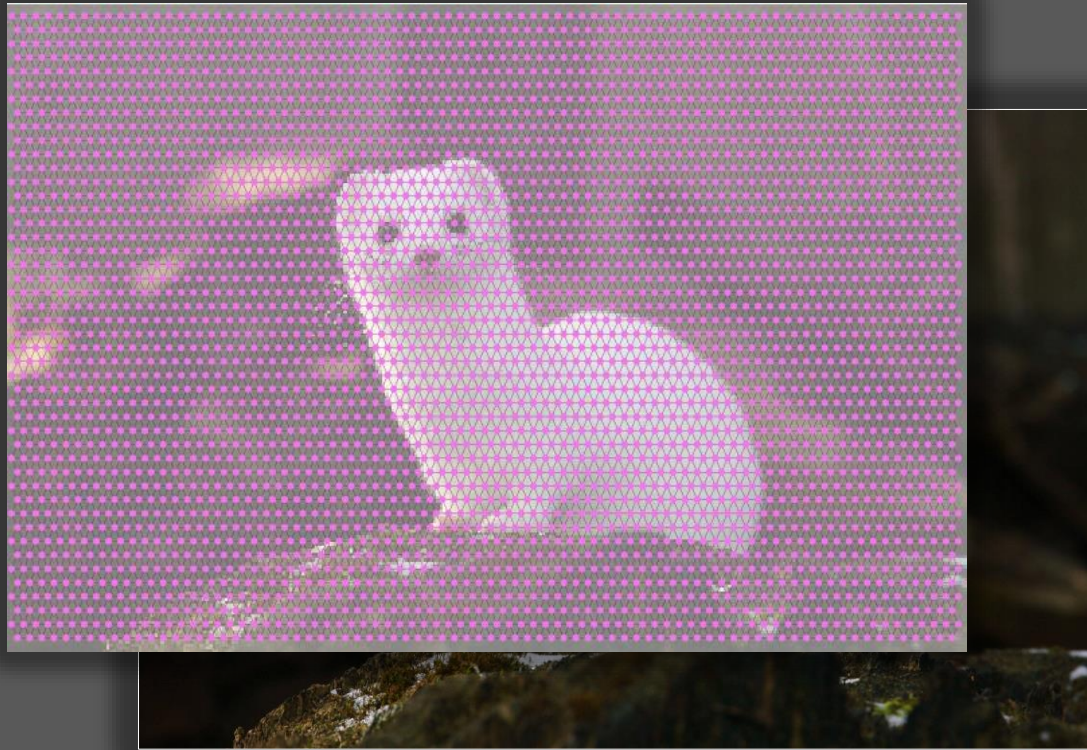




3. run Kruskal's algorithm, and stop when the forest it's building has just a few trees

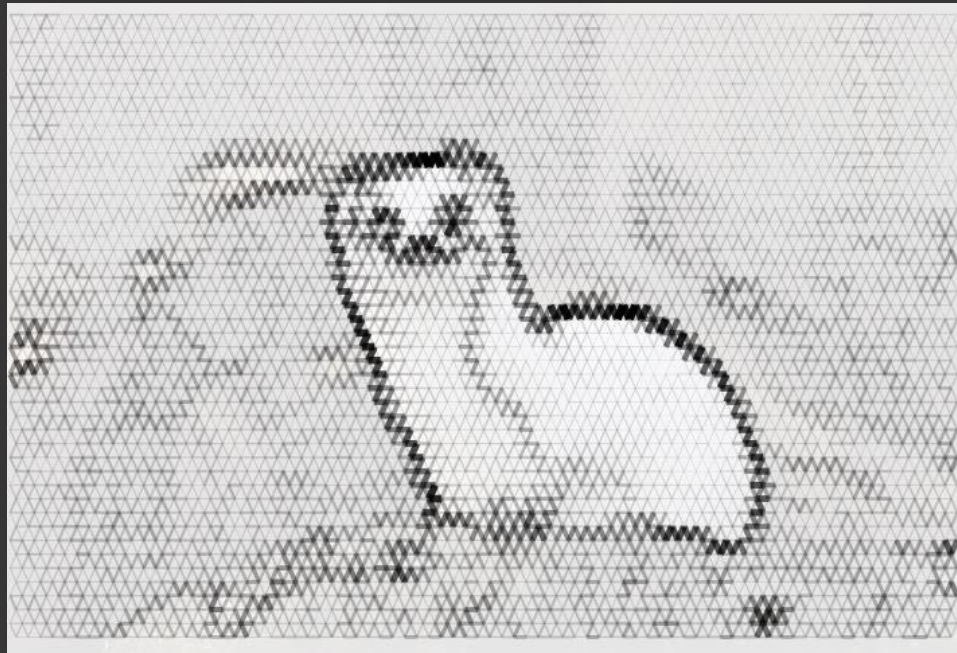


Alternatively ...



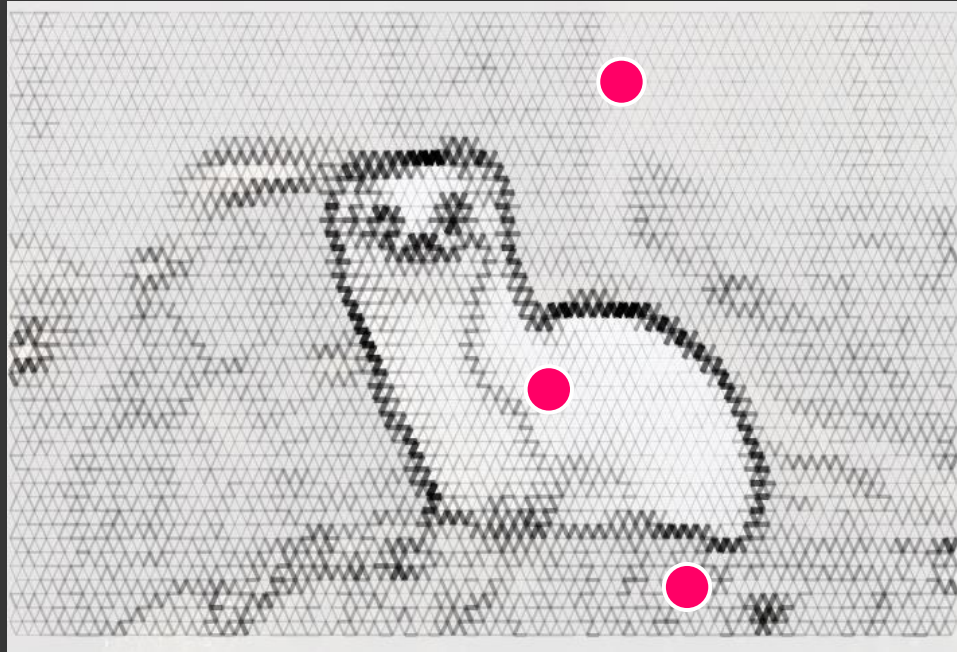
1. define a grid

Alternatively ...



2. measure dissimilarity along edges

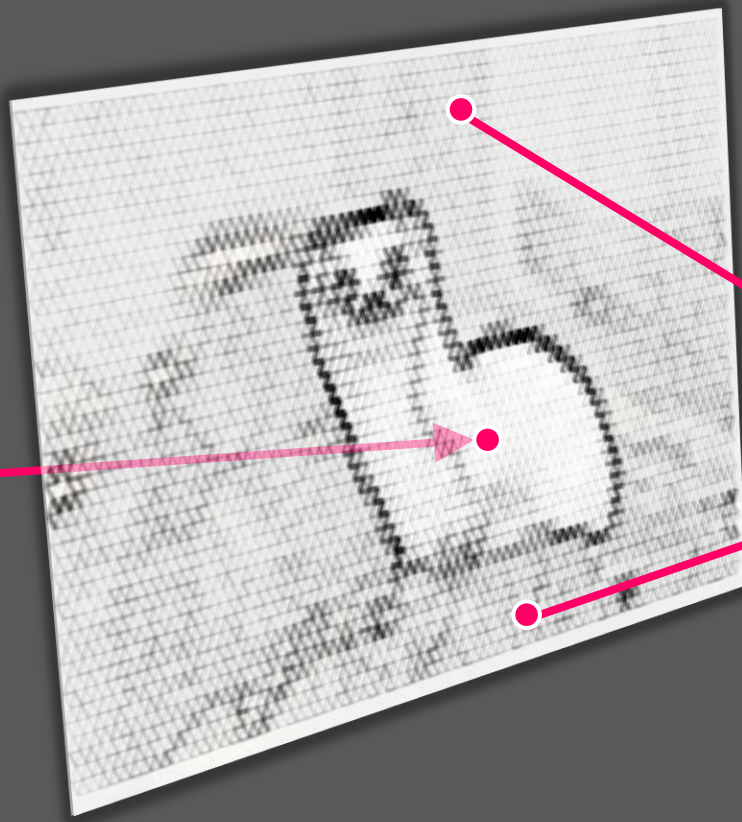
Alternatively ...



3. ask the user to label some “stork” points and some “background” points

Alternatively ...

source



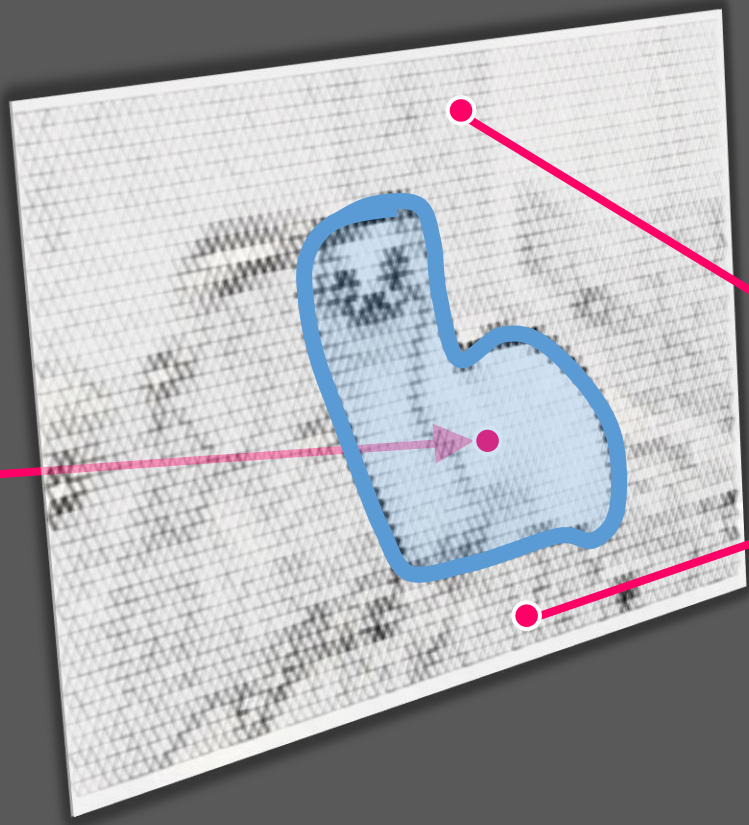
sink



4. set up a flow network

Alternatively ...

source



sink



5. find a minimum-capacity cut