# Minimum Spanning Tree
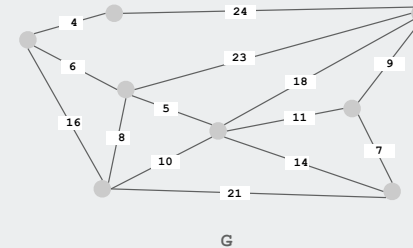
- introduction
- Weighted graph API
- cycles and cuts
- Kruskal's algorithm
- Prim's algorithm
- advanced algorithms
- clustering

---

## Minimum Spanning Tree

**MST.** Given connected graph G with positive edge weights, find a min weight set of edges that connects all of the vertices.



G

---

introduction

weighted graph API
cycles and cuts
Kruskal's algorithm
Prim's algorithm
advanced algorithms
clustering

---

## Minimum Spanning Tree

**MST.** Given connected graph G with positive edge weights, find a min weight set of edges that connects all of the vertices.
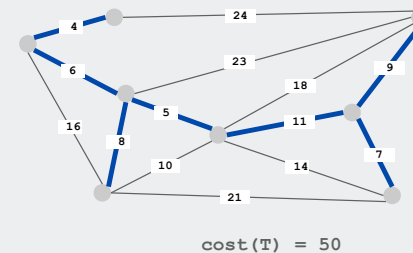


cost(T) = 50

Brute force: Try all possible spanning trees
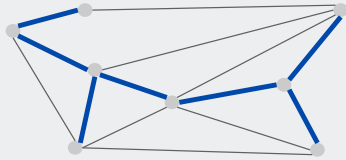- problem 1: not so easy to implement
- problem 2: far too many of them

Ex: [Cayley, 1889]: $V^{V-2}$ spanning trees on the complete graph on V vertices.

## MST Origin

### Otakar Boruvka (1926).
- Electrical Power Company of Western Moravia in Brno.
- Most economical construction of electrical power network.
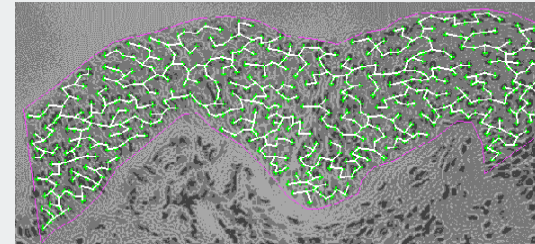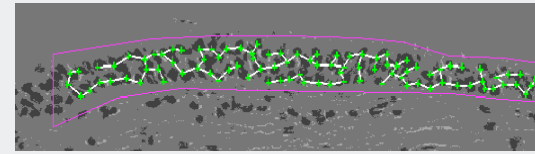- Concrete engineering problem is now a cornerstone problem in combinatorial optimization.



Otakar Boruvka

---

## Medical Image Processing

MST describes arrangement of nuclei in the epithelium for cancer research
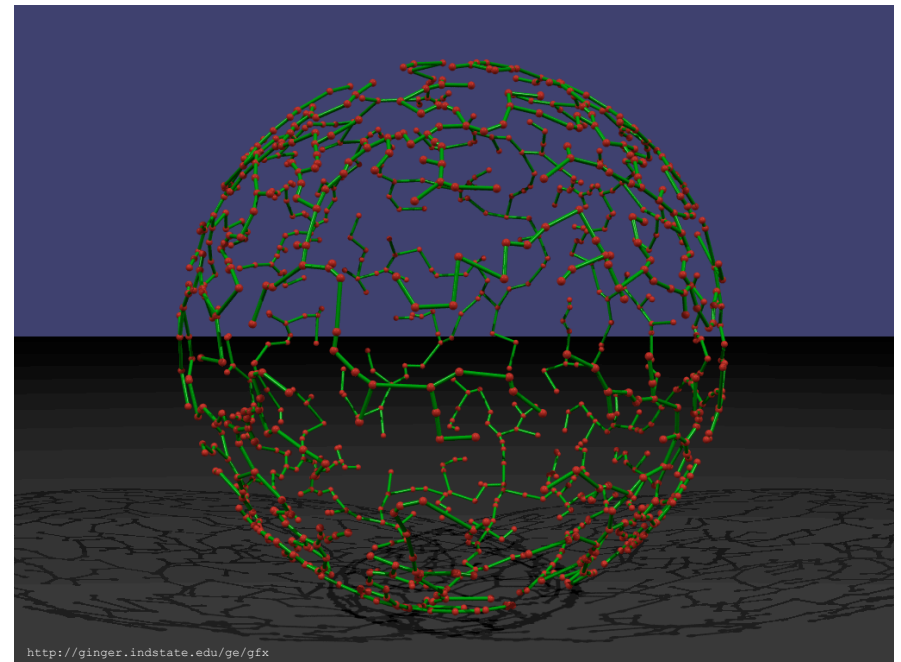


http://www.bccrc.ca/ci/ta01_archlevel.html

---

## Applications

MST is fundamental problem with diverse applications.

- Network design.
  - telephone, electrical, hydraulic, TV cable, computer, road

- Approximation algorithms for NP-hard problems.
  - traveling salesperson problem, Steiner tree

- Indirect applications.
  - max bottleneck paths
  - LDPC codes for error correction
  - image registration with Renyi entropy
  - learning salient features for real-time face verification
  - reducing data storage in sequencing amino acids in a protein
  - model locality of particle interactions in turbulent fluid flows
  - autoconfig protocol for Ethernet bridging to avoid cycles in a network

- Cluster analysis.

---



http://ginger.indstate.edu/ge/gfx

Kruskal's algorithm.  Consider edges in ascending order of cost.
Add the next edge to T unless doing so would create a cycle.

Prim's algorithm.  Start with any vertex s and greedily grow a tree T
from s.  At each step, add the cheapest edge to T that has exactly
one endpoint in T.

Theorem.  Both greedy algorithms compute an MST.

Greed is good.  Greed is right. Greed works.  Greed
clarifies, cuts through, and captures the essence of the
evolutionary spirit."   - Gordon Gecko

---

| public class WeightedGraph | (graph data type) |
|---|---|
| WeightedGraph(int V) | create an empty graph with V vertices |
| void insert(Edge e) | insert edge e |
| Iterable<Edge> adj(int v) | return an iterator over edges incident to v |
| int V() | return the number of vertices |
| String toString() | return a string representation |

```
for (int v = 0; v < G.V(); v++)
{
    for (Edge e : G.adj(v))
    {
        int w = e.other(v);
        // edge v-w
    }
}
```

iterate through all edges (once in each direction)

---

introduction
**weighted graph API**
cycles and cuts
Kruskal's algorithm
Prim's algorithm
advanced algorithms
clustering

---

```
public class Edge implements Comparable<Edge>
{
    public final int v, int w;
    public final double weight;

    public Edge(int v, int w, double weight)
    {
        this.v = v;
        this.w = w;
        this.weight = weight;
    }

    public int either()
    {  return v; }

    public int other(int vertex)
    {
        if (vertex == v) return w;
        else return v;
    }

    public int compareTo(Edge f)
    {
        Edge e = this;
        if      (e.weight < f.weight) return -1;
        else if (e.weight > f.weight) return +1;
        else if (e.weight > f.weight) return  0;
    }
}
```

## Weighted graph: Java implementation

Identical to `Graph.java` but use `Edge` adjacency lists instead of `int`.

```java
public class WeightedGraph
{
   private int V;
   private Sequence<Edge>[] adj;

   public Graph(int V)
   {
      this.V = V;
      adj = (Sequence<Edge>[]) new Sequence[V];
      for (int v = 0; v < V; v++)
         adj[v] = new Sequence<Edge>();
   }

   public void insert(Edge e)
   {
      int v = e.v, w = e.w;
      adj[v].add(e);
      adj[w].add(e);
   }

   public Iterable<Edge> adj(int v)
   {  return adj[v];  }

}
```
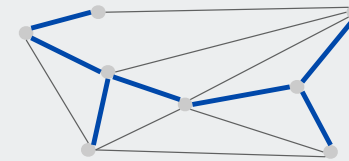
## Spanning Tree

MST.  Given connected graph G with positive edge weights, find a min weight set of edges that connects all of the vertices.

Def.  A spanning tree of a graph G is a subgraph T that is connected and acyclic.
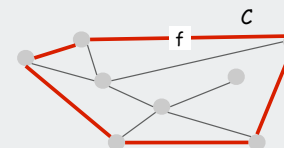


Property.  MST of G is always a spanning tree.

## Greedy Algorithms

Simplifying assumption.  All edge costs $c_e$ are distinct.

Cycle property.  Let C be any cycle, and let f be the max cost edge belonging to C.  Then the MST does not contain f.

Cut property.  Let S be any subset of vertices, and let e be the min cost edge with exactly one endpoint in S.  Then the MST contains e.
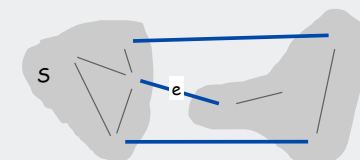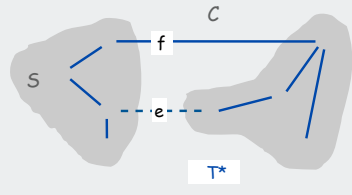


f is not in the MST          e is in the MST

## Cycle Property

Simplifying assumption.  All edge costs $c_e$ are distinct.

Cycle property.  Let C be any cycle, and let f be the max cost edge belonging to C. Then the MST T* does not contain f.

Pf.  [by contradiction]
- Suppose f belongs to T*.  Let's see what happens.
- Deleting f from T* disconnects T*. Let S be one side of the cut.
- Some other edge in C, say e, has exactly one endpoint in S.
- T = T* ∪ { e } – { f } is also a spanning tree.
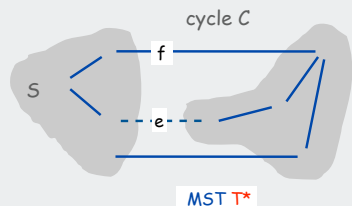- Since $c_e < c_f$, cost(T) < cost(T*).
- Contradicts minimality of T*.  ▪



---

## Cut Property

Simplifying assumption.  All edge costs $c_e$ are distinct.

Cut property.  Let S be any subset of vertices, and let e be the min cost edge with exactly one endpoint in S. Then the MST T* contains e.

Pf.  [by contradiction]
- Suppose e does not belong to T*.  Let's see what happens.
- Adding e to T* creates a (unique) cycle C in T*.
- Some other edge in C, say f, has exactly one endpoint in S.
- T = T* ∪ { e } – { f } is also a spanning tree.
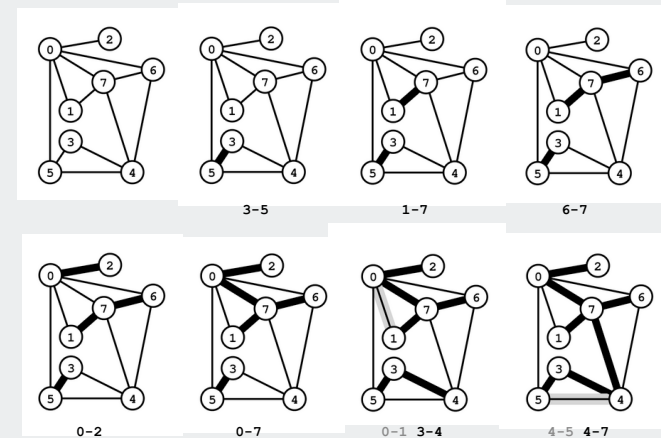- Since $c_e < c_f$, cost(T) < cost(T*).
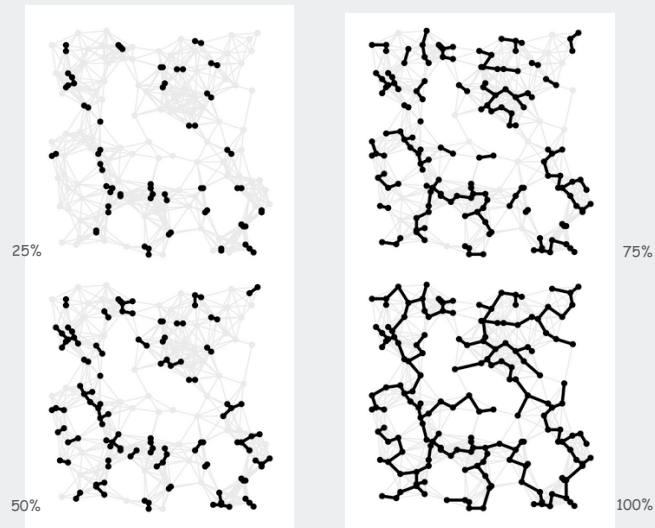- Contradicts minimality of T*.  ▪



---

---

## Kruskal's Algorithm:  Example

Kruskal's algorithm. [Kruskal, 1956]  Consider edges in ascending order of cost. Add the next edge to T unless doing so would create a cycle.

## Kruskal's algorithm example



25%

50%

75%

100%

## Kruskal's algorithm correctness proof

Theorem. Kruskal's algorithm computes the MST.

Pf. [case 2] Suppose that adding e = (v, w) to T does not create a cycle
- let S be the vertices in v's connected component
- w is not in S
- e is the min weight edge with exactly one endpoint in S
- e is in the MST (cut property)

QED

## Kruskal's algorithm correctness proof

Theorem. Kruskal's algorithm computes the MST.

Pf. [case 1] Suppose that adding e to T creates a cycle C
- e is the max weight edge in C (weights come in increasing order)
- e is not in the MST (cycle property)

## Kruskal's algorithm implementation

Q. How to check if adding an edge to T would create a cycle?

A1. Naïve solution: use DFS.
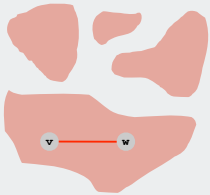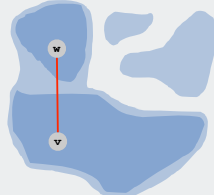- O(V) time per cycle check.
- O(E V) time overall.

Q. How to check if adding an edge to T would create a cycle?

A2. Use the union-find data structure from lecture 1 (!).
- Maintain a set for each connected component.
- If v and w are in same component, then adding v-w creates a cycle.
- To add v-w to T, merge sets containing v and w.



Case 1: adding v-w creates a cycle          Case 2: add v-w to T and merge sets

---

Kruskal running time: Dominated by the cost of the sort.

| Operation | Frequency | Time per op |
|-----------|-----------|-------------|
| sort | 1 | E log E |
| union | V | log* V [†] |
| find | E | log* V [†] |

† amortized bound using weighted quick union with path compression

recall: log* V ≤ 5 in this universe

Remark 1. If edges are already sorted, time is proportional to E log* V

Remark 2. With PQ or quicksort partitioning, time depends on number of edges shorter than longest edge in the MST (may be small in practice).

---

Kruskal's algorithm: Java implementation

```
public class Kruskal
{
    private Sequence<Edge> mst
                   = new Sequence<Edge>();

    public Kruskal(WeightedGraph G)
    {
        Edge[] edges = G.edges();           ← sort edges
        Arrays.sort(edges);

        UnionFind uf = new UnionFind(G.V());
        for (Edge edge : edges)
            if (!uf.find(edge.v, edge.w))   ← greedily add
            {                                  edges to MST
                uf.unite(edge.v, edge.w);
                mst.add(edge);
            }

    }

    public Iterable<Edge> mst()
    {   return mst;   }                      ← return to client iterable
}                                              sequence of edges
```

Easy speedup: Stop as soon as there are V-1 edges in MST.

---

## Prim's Algorithm example

Prim's algorithm. [Jarník 1930, Dijkstra 1957, Prim 1959]
Start with vertex 0 and greedily grow tree T. At each step,
add cheapest edge that has exactly one endpoint in T.



0-2 0-7 0-1 0-6 0-5    0-7 0-1 0-6 0-5    7-1 7-6 7-4 0-5    7-6 7-4 0-5

7-4 0-5    4-3 4-5    3-5

## Prim's Algorithm example



25%    75%

50%    100%

## Prim's algorithm correctness proof

Theorem.  Prim's algorithm computes the MST.
Pf.
- Let S be the subset of vertices in current tree T.
- Prim adds the cheapest edge e with exactly one endpoint in S.
- e is in the MST (cut property)

QED.



S    e

## Prim's algorithm implementation

Q.  How to find cheapest edge with exactly one endpoint in S?

A1.  Brute force:  try all edges.
- O(E) time per spanning tree edge.
- O(E V) time overall.

## Prim's algorithm implementation

Q. How to find cheapest edge with exactly one endpoint in S?

A2. Maintain a priority queue of edges with (at least) one endpoint in S
- Delete min to determine next edge e to add to T.
- Disregard e if both endpoints are in S.
- Upon adding e to T, add to PQ the edges incident to the endpoint not already in S.

Running time.
- log V steps per edge (using a binary heap).
- E log V steps overall.

Note: This is a lazy version of implementation in Algs in Java

lazy:  put all adjacent nodes on PQ
book: first check whether the other endpoint is was put in S
      during the time it was on the PQ

---

## Removing the distinct edge costs assumption

Simplifying assumption.  All edge costs $c_e$ are distinct.
Fact.  Prim and Kruskal don't actually rely on the assumption
       (our proof of correctness does)

Suffices to introduce tie-breaking rule for `compareTo()`.

Approach 1:

```java
public int compareTo(Edge f)
{
    Edge e = this;
    if (e.weight < f.weight) return -1;
    if (e.weight > f.weight) return +1;
    if (e.v < f.v) return -1;
    if (e.v > f.v) return +1;
    if (e.w < f.w) return -1;
    if (e.w > f.w) return +1;
    return  0;
}
```

Approach 2: add tiny random perturbation.

---

## Prim's Algorithm:  Java Implementation

```java
public class LazyPrim
{
    private Sequence<Edge> mst = new Sequence<Edge>();
    public LazyPrim(WeightedGraph G)
    {
        boolean[] marked = new boolean[G.V()];
        MinPQ<Edge> pq = new MinPQ<Edge>();        // marks vertices in s
        int s = 0;                                  // PQ
        marked[s] = true;
        for (Edge e : G.adj(s)) pq.insert(e);       // add to PQ all edges incident to s
        while (!pq.isEmpty())
        {
            Edge e = pq.delMin();
            int v = e.either(), w = e.other(v);
            if (!marked[v])
            {
                mst.add(e); marked(v) = true;
                for (Edge f : G.adj(v)) pq.insert(f);   // add to PQ all edges incident to e.v
            }
            else if (!marked[w])
            {
                mst.add(e); marked(w) = true;
                for (Edge f : G.adj(w)) pq.insert(f);
            }
        }
    }
}
```
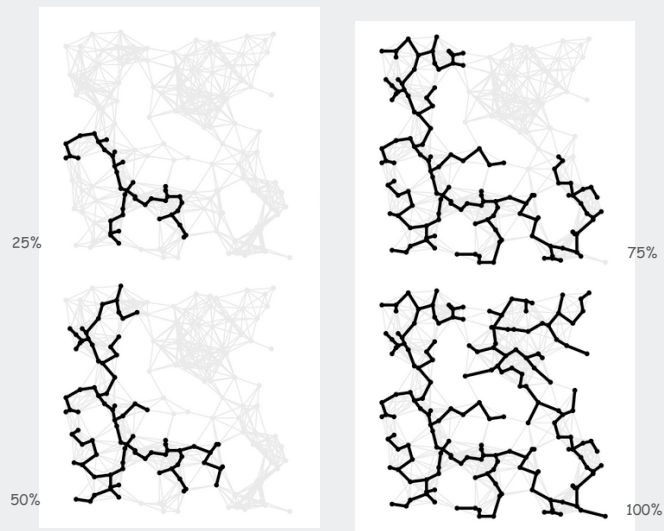
---

introduction
weighted graph API
cycles and cuts
Kruskal's algorithm
Prim's algorithm
**advanced algorithms**
clustering

| Year | Worst Case | Discovered By |
|------|-----------|---------------|
| 1975 | E log log V | Yao |
| 1976 | E log log V | Cheriton-Tarjan |
| 1984 | E log* V,  E + V log V | Fredman-Tarjan |
| 1986 | E log (log* V) | Gabow-Galil-Spencer-Tarjan |
| 1997 | E $\alpha$(V) log $\alpha$(V) | Chazelle |
| 2000 | E $\alpha$(V) | Chazelle |
| 2002 | optimal | Pettie-Ramachandran |
| 20xx | E | ??? |

deterministic comparison based MST algorithms

| Year | Problem | Time | Discovered By |
|------|---------|------|---------------|
| 1976 | Planar MST | E | Cheriton-Tarjan |
| 1992 | MST Verification | E | Dixon-Rauch-Tarjan |
| 1995 | Randomized MST | E | Karger-Klein-Tarjan |

related problems

37

**Key geometric fact.**
Edges of the Euclidean MST are edges of the Delaunay triangulation.

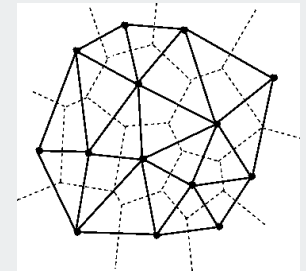**Euclidean MST algorithm.**
- Compute Delaunay triangulation.
- Run Kruskal's MST algorithm on Delaunay edges.

**Running time.**  O(N log N).
- O(N) Delaunay edges since it is planar.
- O(N log N) for Delaunay.
- O(N log N) for Kruskal.



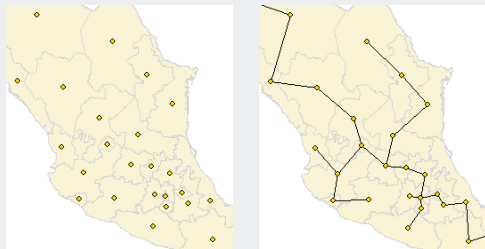**In practice.**  Use any set of edges of size O(N) that contains the Delaunay with high probability

**Lower bound.**  Any comparison-based Euclidean MST algorithm requires $\Omega$(N log N) comparisons.

39

**Euclidean MST.**  Given N points in the plane, find MST connecting them.
- Distances between point pairs are Euclidean distances.



**Brute force.**  Compute  N² / 2  distances and run Prim's algorithm.
**Ingenuity.**  Exploit geometry and do it in O(N log N).

38

introduction
weighted graph API
cycles and cuts
Kruskal's algorithm
Prim's algorithm
advanced algorithms
**clustering**

## Clustering

k-clustering.  Divide a set of objects classify into k coherent groups.
distance function.  numeric value specifying "closeness" of two objects.

Fundamental problem.
 Divide into clusters so that points in different clusters are far apart.



Outbreak of cholera deaths  in London in 1850s.
Reference: Nina Mishra, HP Labs

Applications.
- Routing in mobile ad hoc networks.
- Identify patterns in gene expression.
- Document categorization for web search.
- Similarity searching in medical image databases
- Skycat:  cluster $10^9$ sky objects into stars, quasars, galaxies.

## Single-link clustering algorithm

"Well-known" algorithm for single-link clustering:
- Form V clusters of one object each.
- Find the closest pair of objects such that each object is in a different cluster, and add an edge between them.
- Repeat until there are exactly k clusters.

Observation.  This procedure is precisely Kruskal's algorithm
(stop when there are k connected components).

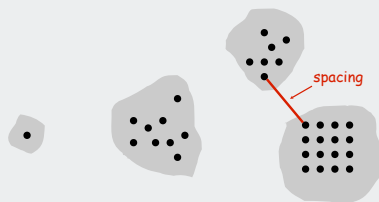Property.  Kruskal's algorithm finds a k-clustering of maximum spacing.

## k-Clustering of maximum spacing

k-clustering.  Divide a set of objects classify into k coherent groups.
distance function.  Numeric value specifying "closeness" of two objects.

Spacing.  Min distance between any pair of points in different clusters.

k-clustering of maximum spacing.
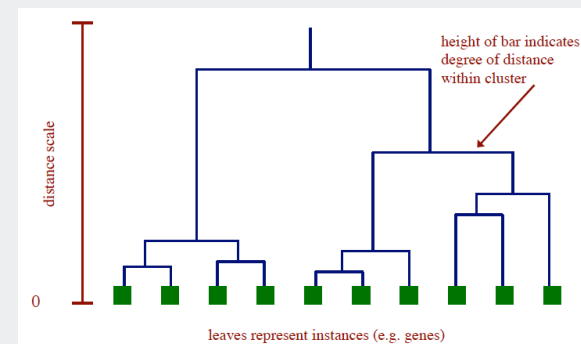Given an integer k, find a k-clustering such that spacing is maximized.



spacing

k = 4

## Dendrogram

Dendrogram.
Scientific visualization of hypothetical sequence of evolutionary events.

- Leaves = genes.
- Internal nodes = hypothetical ancestors.



height of bar indicates degree of distance within cluster

distance scale

0

leaves represent instances (e.g. genes)

Reference:  http://www.biostat.wisc.edu/bmi576/fall-2003/lecture13.pdf

# Dendrogram of cancers in human

Tumors in similar tissues cluster together.



Gene 1

Gene n

Skin  Liver  Lung  Breast Tumors  Breast  Normal  Kidney  Prostate  Brain  APL  Ovary
                    Luminal        Tumors  Breast
                                   Basal

Reference: Botstein & Brown group

☐ gene expressed
■ gene not expressed

45