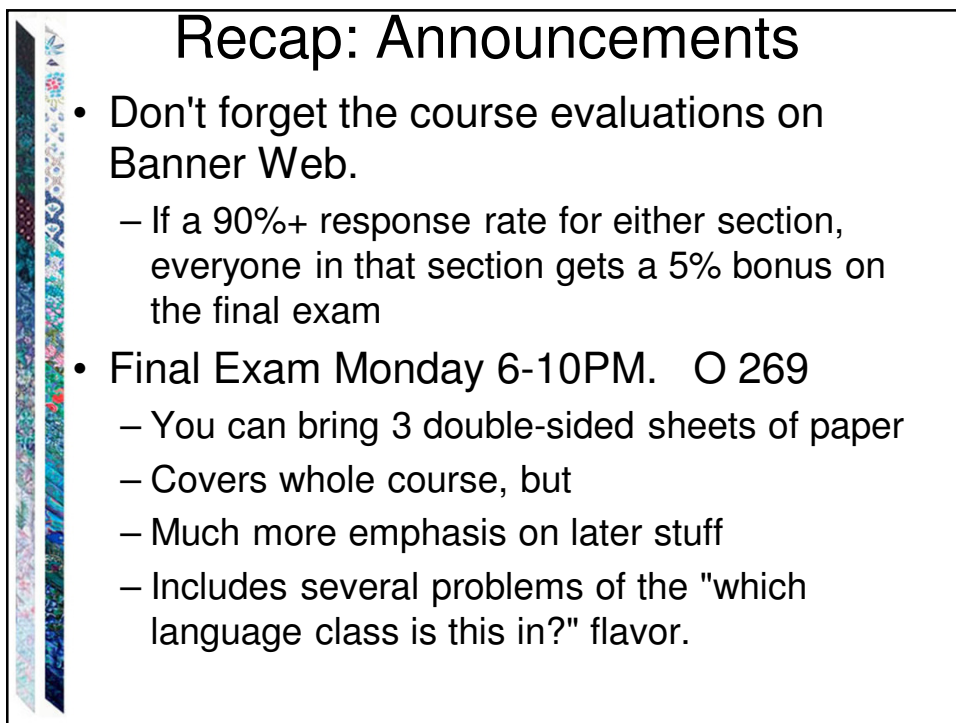


MA/CSSE 474

Theory of Computation

Computational Complexity Continued
P and NP

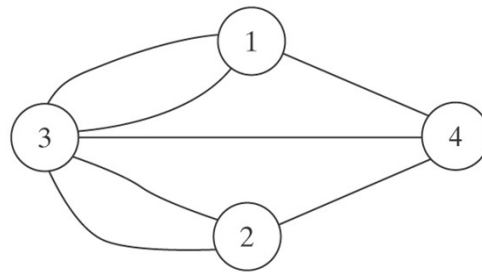


Recap: Announcements

- Don't forget the course evaluations on Banner Web.
 - If a 90%+ response rate for either section, everyone in that section gets a 5% bonus on the final exam
- Final Exam Monday 6-10PM. O 269
 - You can bring 3 double-sided sheets of paper
 - Covers whole course, but
 - Much more emphasis on later stuff
 - Includes several problems of the "which language class is this in?" flavor.

Graph Languages

- $\text{CONNECTED} = \{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is connected} \}$.
- $\text{HAMILTONIANCIRCUIT} = \{ \langle G \rangle : G \text{ is an undirected graph that contains a *Hamiltonian circuit*} \}$.



Characterizing Optimization Problems as Languages

- $\text{TSP-DECIDE} = \{ \langle G, cost \rangle : \langle G \rangle \text{ encodes an undirected graph with a positive distance attached to each of its edges and } G \text{ contains a Hamiltonian circuit whose total cost is less than } \langle cost \rangle \}$.

Choosing A Model of Computation

We'll use Turing machines:

- Tape alphabet size?
- How many tapes?
- Deterministic vs. nondeterministic?

Measuring Time and Space Requirements

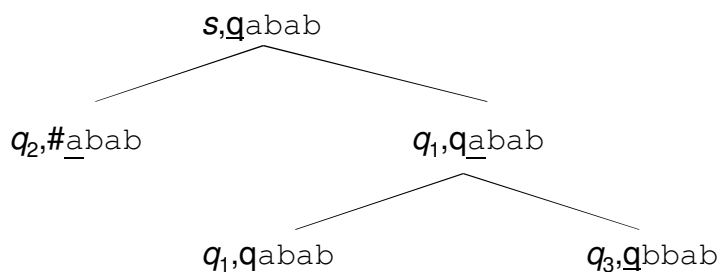
$timereq(M)$ is a function of n :

- If M is a **deterministic** Turing machine that halts on all inputs, then:

$timereq(M) = f(n)$ = the maximum number of steps that M executes on any input of length n .

Measuring Time and Space Requirements

- If M is a **nondeterministic** Turing machine all of whose computational paths halt on all inputs, then:



$timereq(M) = f(n)$ = the number of steps on the longest path that M executes on any input of length n .

Measuring Time and Space Requirements

$spacereq(M)$ is a function of n :

- If M is a **deterministic** Turing machine that halts on all inputs, then:

$spacereq(M) = f(n)$ = the maximum number of different tape squares that M reads on any input of length n .

- If M is a **nondeterministic** Turing machine all of whose computational paths halt on all inputs, then:

$spacereq(M) = f(n)$ = the maximum number of different tape squares that M reads on any path that it executes on any input of length n .

Algorithmic Gaps

We'd like to show for a language L :

1. Upper bound: There exists an algorithm that decides L and that has complexity C_1 .
2. Lower bound: Any algorithm that decides L must have complexity at least C_2 .
3. $C_1 = C_2$.

If $C_1 = C_2$, we are done. Often, we're not done.

Example: Sorting (SDT, merge, heap, sleep)

Algorithmic Gaps

Example: TSP

- Upper bound: $timereq \in \mathcal{O}(2^{(n^k)})$.
- Don't have a lower bound that says polynomial isn't possible.

We group languages by what we know. And then we ask:
"Is class CL_1 equal to class CL_2 ?"

A Simple Example of Polynomial Speedup

Given a list of n numbers, find the minimum and the maximum elements in the list.

Or, as a language recognition problem:

$L = \{ \langle \text{list of numbers}; \text{number}_1; \text{number}_2 \rangle : \text{number}_1 \text{ is the minimum element of the list and } \text{number}_2 \text{ is the maximum element} \}$.

$(23, 45, 73, 12, 45, 197; 12; 197) \in L$.

A Simple Example of Polynomial Speedup

The straightforward approach:

$\text{simplecompare}(\text{list: list of numbers}) =$

$\text{max} = \text{list}[1].$

$\text{min} = \text{list}[1].$

For $i = 2$ to $\text{length}(\text{list})$ do:

 If $\text{list}[i] < \text{min}$ then $\text{min} = \text{list}[i].$

 If $\text{list}[i] > \text{max}$ then $\text{max} = \text{list}[i].$

Requires $2(n-1)$ comparisons. So simplecompare is $\mathcal{O}(n)$.

But we can solve this problem in $(3/2)(n-1)$ comparisons.

How?

A Simple Example of Polynomial Speedup

efficientcompare(list: list of numbers) =

max = list[1].

min = list[1].

For $i = 3$ to $\text{length}(\text{list})$ by 2 do:

 If $\text{list}[i] < \text{list}[i-1]$ then:

 If $\text{list}[i] < \text{min}$ then $\text{min} = \text{list}[i]$.

 If $\text{list}[i-1] > \text{max}$ then $\text{max} = \text{list}[i-1]$.

 Else:

 If $\text{list}[i-1] < \text{min}$ then $\text{min} = \text{list}[i-1]$.

 If $\text{list}[i] > \text{max}$ then $\text{max} = \text{list}[i]$.

If $\text{length}(\text{list})$ is even then check the last element.

Requires $3/2(n-1)$ comparisons.

String Search

t: a b c a b a b c a b d

p: a b c d

 a b c d

 a b c d

 . . .

String Search

```

simple-string-search( $t, p$ : strings) =
   $i = 0$ .
   $j = 0$ .
  While  $i \leq |t| - |p|$  do:
    While  $j < |p|$  do:
      If  $t[i+j] = p[j]$  then  $j = j + 1$ .
      Else exit this loop.
      If  $j = |p|$  then halt and accept.
      Else:
         $i = i + 1$ .
         $j = 0$ .
  Halt and reject.

```

Let n be $|t|$ and let m be $|p|$. In the worst case (in which it doesn't find an early match), *simple-string-search* will go through its outer loop almost n times and, for each of those iterations, it will go through its inner loop m times.

So $\text{timereq}(\text{simple-string-search}) \in \mathcal{O}(nm)$. K-M-P algorithm is $\mathcal{O}(n+m)$

Replacing an Exponential Algorithm with a Polynomial One

- Context-free parsing can be done in $\mathcal{O}(n^3)$ time instead of $\mathcal{O}(2^n)$ time. (CYK algorithm)
- Finding the greatest common divisor of two integers can be done in $\mathcal{O}(\log_2(\max(n, m)))$ time instead of exponential time.

The Language Class P

$L \in P$ iff

- there exists some deterministic Turing machine M that decides L , and
- $\text{timereq}(M) \in \mathcal{O}(n^k)$ for some k .

We'll say that L is **tractable** iff it is in P.

Closure under Complement

Theorem: The class P is closed under complement.

Proof: If M accepts L in polynomial time, swap accepting and non accepting states to accept $\neg L$ in polynomial time.

Defining Complement

- $\text{CONNECTED} = \{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is connected} \}$ is in P.
- $\text{NOTCONNECTED} = \{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is not connected} \}$.
- $\neg\text{CONNECTED} = \text{NOTCONNECTED} \cup \{ \text{strings that are not syntactically legal descriptions of undirected graphs} \}$.

$\neg\text{CONNECTED}$ is in P by the closure theorem. What about NOTCONNECTED?

If we can check for legal syntax in polynomial time, then we can consider the universe of strings whose syntax is legal. Then we can conclude that NOTCONNECTED is in P if CONNECTED is.

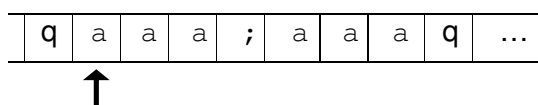
Languages That Are in P

- Every regular language.
- Every context-free language since there exist context-free parsing algorithms that run in $\mathcal{O}(n^3)$ time.
- Others:
 - $A^n B^n C^n$
 - Nim

To Show That a Language Is In P

- Describe a one-tape, deterministic Turing machine.
- It may use multiple tapes. Price:
- State an algorithm that runs on a conventional computer.
Price:

How long does it take to compare two strings?



Bottom line: If ignoring polynomial factors, then just describe a deterministic algorithm.

Regular Languages

Theorem: Every regular language can be decided in linear time. So every regular language is in P.

Proof: If L is regular, there exists some DFSA M that decides it. Construct a deterministic TM M' that simulates M , moving its read/write head one square to the right at each step. When M' reads a q , it halts. If it is in an accepting state, it accepts; otherwise it rejects.

On any input of length n , M' will execute $n + 2$ steps.

So $\text{time}_{req}(M') \in \mathcal{O}(n)$.

Context-Free Languages

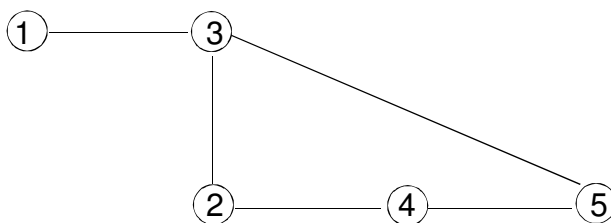
Theorem: Every context-free language can be decided in $\mathcal{O}(n^{18})$ time. So every context-free language is in P.

Proof: The Cocke-Kasami-Younger (CKY) algorithm can parse any context-free language in time that is $\mathcal{O}(n^3)$ if we count operations on a conventional computer. That algorithm can be simulated on a standard, one-tape Turing machine in $\mathcal{O}(n^{18})$ steps.

WE could get bogged down in the details of this, but we won't!

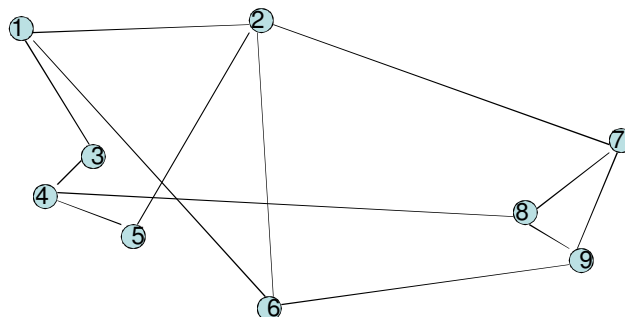
Graph Languages

Represent a graph $G = (V, E)$ as a list of edges:



101/1/11/11/10/10/100/100/101/11/101

Graph Languages



CONNECTED =
 $\{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is connected} \}$.

Is CONNECTED in P?

CONNECTED is in P

$connected(\langle G = (V, E) \rangle =$

1. Set all vertices to be unmarked.
2. Mark vertex 1.
3. Initialize L to $\{1\}$.
4. Initialize *marked-vertices-counter* to 1.
5. Until L is empty do:
 - 5.1. Remove the first element from L .
Call it *current-vertex*.
 - 5.2. For each edge e
that has *current-vertex* as an endpoint do:
Call the other endpoint of e *next-vertex*.
If *next-vertex* is not already marked then do:
Mark *next-vertex*.
Add *next-vertex* to L .
Increment *marked-vertices-counter* by 1.
6. If *marked-vertices-counter* = $|V|$ accept. Else reject.

Analyzing *connected*

- Step 1 takes time that is $\mathcal{O}(|V|)$.
- Steps 2, 3, and 4 each take constant time.
- The loop of step 5 can be executed at most $|V|$ times.
 - Step 5.1 takes constant time.
 - Step 5.2 can be executed at most $|E|$ times. Each time, it requires at most $\mathcal{O}(|V|)$ time.
- Step 6 takes constant time.

So *timereq(connected)* is:

$$|V| \cdot \mathcal{O}(|E|) \cdot \mathcal{O}(|V|) = \mathcal{O}(|V|^2|E|).$$

But $|E| \leq |V|^2$.

So *timereq(connected)* is:
 $\mathcal{O}(|V|^4)$.

```

connected(<G = (V, E)>) =
1. Set all vertices to be unmarked.
2. Mark vertex 1.
3. Initialize L to {1}.
4. Initialize marked-vertices-counter to 1.
5. Until L is empty do:
  5.1. Remove the first element from L.
      Call it current-vertex.
  5.2. For each edge e
      that has current-vertex as an endpoint do:
        Call the other endpoint of e next-vertex.
        If next-vertex is not already marked then do:
          Mark next-vertex.
          Add next-vertex to L.
          Increment marked-vertices-counter by 1
6. If marked-vertices-counter = |V| accept. Else reject.
  
```

Primality Testing

RELATIVELY-PRIME =

$\{ \langle n, m \rangle : n \text{ and } m \text{ are integers that are relatively prime} \}$.

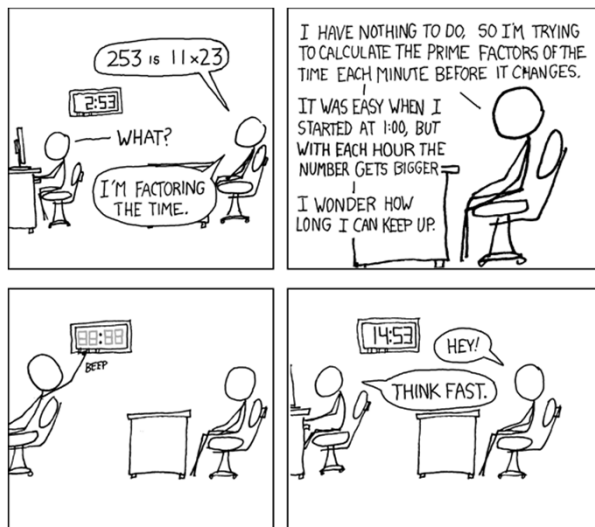
PRIMES =

$\{ w : w \text{ is the binary encoding of a prime number} \}$

COMPOSITES =

$\{ w : w \text{ is the binary encoding of a nonprime number} \}$

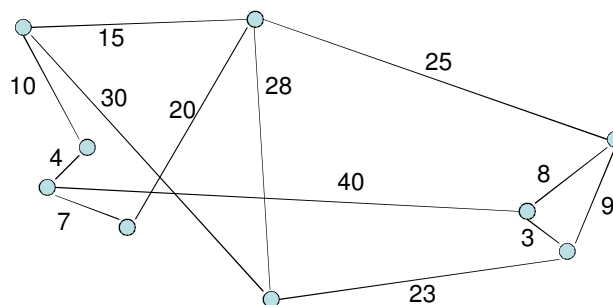
But Finding Factors Remains Hard



<http://xkcd.com/247/>

Returning to TSP

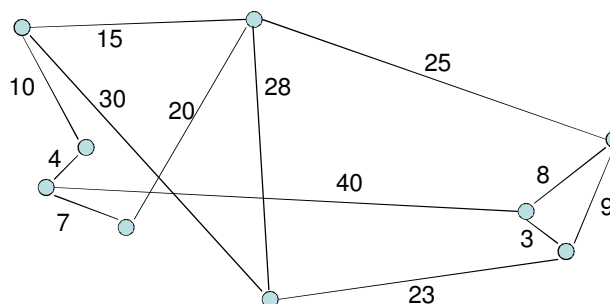
TSP-DECIDE = $\{ \langle G, cost \rangle : \langle G \rangle$ encodes an undirected graph with a positive distance attached to each of its edges and G contains a Hamiltonian circuit whose total cost is less than $\langle cost \rangle$.



An NDTM to decide TSP-DECIDE:

Returning to TSP

An NDTM to decide TSP-DECIDE:



1. For $i = 1$ to $|V|$ do:
 - Choose a vertex that hasn't yet been chosen.
2. Check that the path defined by the chosen sequence of vertices is a Hamiltonian circuit through G with distance less than cost.

TSP and Other Problems Like It

TSP-DECIDE, and other problems like it, share three properties:

1. The problem can be solved by searching through a space of partial solutions (such as routes). The size of this space grows exponentially with the size of the problem.
2. No better technique for finding an exact solution is known.
3. But, if a proposed solution were suddenly to appear, it could be checked for correctness very efficiently.

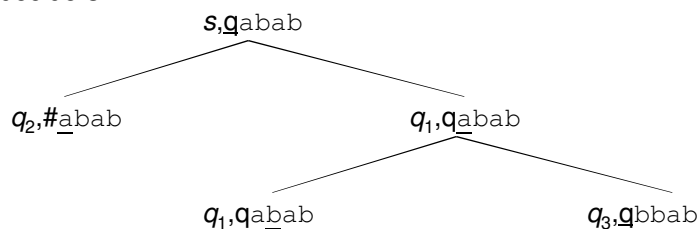
The Language Class NP

Nondeterministic deciding:

$L \in \text{NP}$ iff:

- there is some NDTM M that decides L , and
- $\text{timereq}(M) \in \mathcal{O}(n^k)$ for some k .

NDTM deciders:



TSP Again

$\text{TSP-DECIDE} = \{ \langle G, \text{cost} \rangle : \langle G \rangle \text{ encodes an undirected graph with a positive distance attached to each of its edges and } G \text{ contains a Hamiltonian circuit whose total cost is less than } \langle \text{cost} \rangle \}$.

Suppose some Oracle presented a candidate path c :

$\langle G, \text{cost}, v_1, v_7, v_4, v_3, v_8, v_5, v_2, v_6, v_1 \rangle$

How long would it take to verify that c proves that:

$\langle G, \text{cost} \rangle$ is in TSP-DECIDE?

Deterministic Verifying

A Turing machine V is a **verifier** for a language L iff:

$w \in L$ iff $\exists c \langle w, c \rangle \in L(V)$.

We'll call c a **certificate**.

Deterministic Verifying

An alternative definition for the class NP:

$L \in \text{NP}$ iff there exists a deterministic TM V such that:

- V is a verifier for L , and
- $\text{timereq}(V) \in \mathcal{O}(n^k)$ for some k .

ND Deciding and D Verifying

Theorem: These two definitions are equivalent:

- (1) $L \in \text{NP}$ iff there exists a nondeterministic, polynomial-time TM that decides it.
- (2) $L \in \text{NP}$ iff there exists a deterministic, polynomial-time verifier for it.

Proof: We skip it

Proving That a Language is in NP

- Exhibit an NDTM to decide it.
- Exhibit a DTM to verify it.

Example

- $SAT = \{w : w \text{ is a Boolean wff and } w \text{ is satisfiable}\}$ is in NP.

$$F_1 = P \wedge Q \wedge \neg R ?$$

$$F_2 = P \wedge Q \wedge R ?$$

$$F_3 = P \wedge \neg P ?$$

$$F_4 = P \wedge (Q \vee \neg R) \wedge \neg Q ?$$

$$SAT-decide(F_4) =$$

$$SAT-verify (\langle F_4, (P = True, Q = False, R = False) \rangle) =$$

3-SAT

- A **literal** is either a variable or a variable preceded by a single negation symbol.
- A **clause** is either a single literal or the disjunction of two or more literals.
- A wff is in **conjunctive normal form** (or CNF) iff it is either a single clause or the conjunction of two or more clauses.
- A wff is in **3-conjunctive normal form** (or 3-CNF) iff it is in conjunctive normal form and each clause contains exactly three literals.

3-SAT

| | 3-CNF | CNF |
|---|-------|-----|
| $(P \vee \neg Q \vee R)$ | • | • |
| $(P \vee \neg Q \vee R) \wedge (\neg P \vee Q \vee \neg R)$ | • | • |
| P | | • |
| $(P \vee \neg Q \vee R \vee S) \wedge (\neg P \vee \neg R)$ | | • |
| $P \rightarrow Q$ | | |
| $(P \wedge \neg Q \wedge R \wedge S) \vee (\neg P \wedge \neg R)$ | | |
| $\neg(P \vee Q \vee R)$ | | |

Every wff can be converted to an equivalent wff in CNF.

- 3-SAT = { w : w is a wff in Boolean logic, w is in 3-conjunctive normal form, and w is satisfiable }.

3-SAT is in NP

The Relationship Between P and NP

Is $P = NP$?

Here are some things we know:

$P \subseteq NP \subseteq PSPACE \subseteq EXPTIME$

$P \neq EXPTIME$

[The Millenium Prize](#)

Using Reduction in Complexity Proofs

A **mapping reduction** R from L_1 to L_2 is a

- Turing machine that
- implements some computable function f with the property that:

$$\forall x (x \in L_1 \leftrightarrow f(x) \in L_2).$$

If $L_1 \leq L_2$ and M decides L_2 , then:

$$C(x) = M(R(x)) \text{ will decide } L_1.$$

Using Reduction in Complexity Proofs

If R is deterministic polynomial time then:

$$L_1 \leq_P L_2.$$

And, whenever such an R exists:

- L_1 must be in P if L_2 is: if L_2 is in P then there exists some deterministic, polynomial-time Turing machine M that decides it. So $M(R(x))$ is also a deterministic, polynomial-time Turing machine and it decides L_1 .
- L_1 must be in NP if L_2 is: if L_2 is in NP then there exists some nondeterministic, polynomial-time Turing machine M that decides it. So $M(R(x))$ is also a nondeterministic, polynomial-time Turing machine and it decides L_1 .

Why Use Reduction?

Given $L_1 \leq_P L_2$, we can use reduction to:

- Prove that L_1 is in P or in NP because we **already know** that L_2 is.
- Prove that L_1 would be in P or in NP if we **could somehow show** that L_2 is. When we do this, we cluster languages of similar complexity (even if we're not yet sure what that complexity is). In other words, L_1 is no harder than L_2 is.

NP-Completeness

A language L might have these properties:

1. L is in NP.
 2. Every language in NP is deterministic, polynomial-time reducible to L .
- L is **NP-hard** iff it possesses property 2.
An NP-hard language is at least as hard as any other language in NP.
 - L is **NP-complete** iff it possesses *both* property 1 and property 2.
All NP-complete languages can be viewed as being equivalently hard.

NP-Complete Languages

- SUBSET-SUM = $\{\langle S, k \rangle : S \text{ is a multiset of integers, } k \text{ is an integer, and there exists some subset of } S \text{ whose elements sum to } k\}$.
- SET-PARTITION = $\{\langle S \rangle : S \text{ is a multiset of objects each of which has an associated cost and there exists a way to divide } S \text{ into two subsets, } A \text{ and } S - A, \text{ such that the sum of the costs of the elements in } A \text{ equals the sum of the costs of the elements in } S - A\}$.
- KNAPSACK = $\{\langle S, v, c \rangle : S \text{ is a set of objects each of which has an associated cost and an associated value, } v \text{ and } c \text{ are integers, and there exists some way of choosing elements of } S \text{ (duplicates allowed) such that the total cost of the chosen objects is at most } c \text{ and their total value is at least } v\}$.

NP-Complete Languages

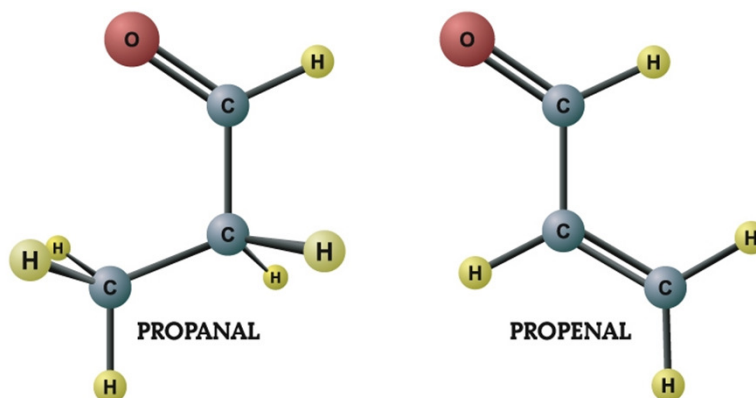
- TSP-DECIDE.
- HAMILTONIAN-PATH = $\{\langle G \rangle : G \text{ is an undirected graph and } G \text{ contains a Hamiltonian path}\}$.
- HAMILTONIAN-CIRCUIT = $\{\langle G \rangle : G \text{ is an undirected graph and } G \text{ contains a Hamiltonian circuit}\}$.
- CLIQUE = $\{\langle G, k \rangle : G \text{ is an undirected graph with vertices } V \text{ and edges } E, k \text{ is an integer, } 1 \leq k \leq |V|, \text{ and } G \text{ contains a } k\text{-clique}\}$.
- INDEPENDENT-SET = $\{\langle G, k \rangle : G \text{ is an undirected graph and } G \text{ contains an independent set of at least } k \text{ vertices}\}$.

NP-Complete Languages

- SUBGRAPH-ISOMORPHISM = $\{ \langle G_1, G_2 \rangle : G_1 \text{ is isomorphic to some subgraph of } G_2 \}$.

Two graphs G and H are **isomorphic** to each other iff there exists a way to rename the vertices of G so that the result is equal to H . Another way to think about isomorphism is that two graphs are isomorphic iff their drawings are identical except for the labels on the vertices.

SUBGRAPH-ISOMORPHISM



NP-Complete Languages

- BIN-PACKING = $\{ \langle S, c, k \rangle : S \text{ is a set of objects each of which has an associated size and it is possible to divide the objects so that they fit into } k \text{ bins, each of which has size } c \}$.

BIN-PACKING

In three dimensions:



Proving that L is NP-Complete

NPL_1 NPL_2 NPL_3 NPL_4 $NPL\dots$

L_1



L_2

Theorem:

If:

L_1 is NP-complete,

$L_1 \leq_P L_2$, and

L_2 is in NP,

Then L_2 is also NP-complete.