

Seminar

Gems of Theoretical Computer Science

LOGSPACE, RandomWalks on Graphs and Universal
Traversal Sequences

Summer term 2014

Institute of Theoretical Informatics (ITI)
Department of Informatics
Karlsruhe Institute of Technology (KIT)

Author: Patrick Niklaus
Betreuer: Jun.-Prof. Dr. Henning Meyerhenke

Copyright © 2014 ITI and Patrick Niklaus

Institute of Theoretical Informatics (ITI)
Department of Informatics
Karlsruhe Institute of Technology
Am Fasanengarten 5
76128 Karlsruhe

Contents

1	LOGSPACE, RandomWalks on Graphs and Universal Traversal Sequences	1
	(Patrick Niklaus)	
2	Time and space complexity	3
2.1	Turing machine models	3
2.2	Time and space bounds	4
2.3	Decision problems and Complexity classes	4
2.4	$NL \subseteq P$	5
2.5	Reducibility and NL-completeness	5
3	PATH and UPATH	7
3.1	PATH	7
3.1.1	NL-completeness	7
3.2	UPATH	8
3.3	UPATH vs. PATH	8
4	RL	9
4.1	Randomized vs. Non-Deterministic	9
4.2	Proof that poly-running time bound is required	9
5	Random Walk	11
5.1	The algorithm	11
5.2	Application on PATH	11
5.3	Correctness for UPATH	12
6	Universal Traversal Sequence	15
6.1	Definition	15
6.2	Relation to previous proof	15
	References	15

Chapter 1

LOGSPACE, RandomWalks on Graphs and Universal Traversal Sequences

Patrick Niklaus

In this paper I will present the collected findings of Aleliunas et. al. [2] on the space complexity of deciding whether there exists a path between two nodes a, b in graphs.

Interestingly the complexity of this decision problem is different for directed and undirected graphs. These differences can be exploited to characterize the two space-bounded complexity classes L and NL .

In the case of an undirected graph, we can employ an algorithm called *random walk* which does just that. Unlike common path finding algorithms like the famous Dijkstra algorithm, this algorithm only needs a logarithmic amount of storage to compute, and is inherently simplistic in its description: Instead of choosing the next neighbour node to explore based on certain criteria, we just pick one at random and check if we reached our destination. It is not intuitively clear that this algorithm can decide whether there is a path between a and b correctly. Infact it can only do so with a certain probability, it is an *randomized* algorithm. However, by choosing a sufficient number of steps we can make sure it decides correctly with a probability higher than $\frac{1}{2}$, which we can still increase by just running the algorithm multiple times.

Furthermore this is deeply related to finding “routing directions”, that visit all nodes in a graph and work in all graphs with a certain size, so called *universal traversal sequences*. We will show that for every graph that is d – *connected*, there exists such a sequences and its length is only polynomial in the number of nodes.

Further Research A lot of research has be done by Aldous et. al on providing tighter lower bounds for the minium number of steps it takes to visit all nodes in graph, both for special graphs e.g. b-ary trees [1] and general undirected graphs [4]. Also are several attempts to get a upper bound on the space-complexity of *UPATH* using only *deterministic* TMs, most recently Armoni et al [3] gave a space bound of $O((\log n)^{\frac{4}{3}})$.

Applications Apart from the theoretical implications for complexity theory, RandomWalk and Universal Traversal Sequences

have well known applications in Artificial Intelligence for state space exploration [5].

Chapter 2

Time and space complexity

2.1 Turing machine models

For this paper we use the common formal definition of a *turing machine*, which we will extent to have 3 seperate tapes: A read-only reading tape, a work tape that can be read and written, and a write-only output tape. We do this, to be able to ignore any reading operations done on the input, which is important to formalize sub-linear space bounds. We also require the reading head of the input tape to stay only on the non-blank part of the input, which is important to get a bound for the configurations of such a turing machine.

A *turing machine* is a tuple $M = (Q, \Gamma, b, \Sigma, q_0, F, \delta)$ with:

Q Set of states

Γ Set of tape symbols

$b \in \Gamma$ Blank symbol

$\Sigma \subseteq \Gamma$ Set of input symbols

$q_0 \in Q$ Initial state

$F \subseteq Q$ Set of halting states.

$\delta : \Gamma^2 \times Q \longrightarrow \Gamma^2 \times Q \times \{L, R, N\}^3$ (partial) transition function, note that this function can read and write two characters in each step, because we can read from two tapes and write to two tapes.

We will also use the notion of *turing machine acceptors*, which we can formalize by choosing a subset of the halting states as *accepting* states.

As you might have noticed, this TM definition is that of a deterministic TM. For non-deterministic TM we simply remove the restriction that δ is a (partial) function. Instead δ is a general relation, which for example implies that there can be more than one next configuration.

This results in the property that a computation of a non-deterministic TM is not a single sequence of configurations, but rather a set of sequences. A deterministic TM accepts if it halts in an accepting state. To redefine this for a NTM, we say it accepts an input, if there *exists at least one* sequence of configurations that ends in an accepting state and the NTM halts.

2.2 Time and space bounds

Before we continue, it is important to establish the notion of *running time* and *space usage* since we will use it later to differentiate certain classes of problems.

For a *deterministic* TM it is easy enough to define running time as the number of steps (i.e. the number of configurations) it takes before it halts. For a *non-deterministic* TM we need to expand this to the *maximum* number of configurations that lead to an accepting state, which is equivalent to the depth of the computation tree.

Space usage in a 3 tape model is defined as the number cells on the *working tape* that were visited by the reading head during the computation. As before, we use the maximum of all possible halting computations in the non-deterministic case.

2.3 Decision problems and Complexity classes

A common class of problems are so called *Decision Problems*, that describe (in the most general case) the following problems: Given a language $L \subseteq \Sigma^*$ and a word $w \in \Sigma^*$: Is $w \in L$? To formalize this, we can use *turing machine acceptors*, that halt in an accepting state if $w \in L$ and in a non-accepting state otherwise.

Since we have defined a notion of running time and space usage for both turing machine models, we can use this to classify decision problems based on running time and space usage. In general we say, a language L is in a complexity class A if and only if the corresponding decision problem satisfies a certain constraint. For this paper, we are mainly interested in four complexity classes (note when we say *polynomially-bounded* we mean: bounded by a polynomial in the input length n):

- P** All decision problems that can be solved by a *deterministic* turing machine using only polynomially-bounded many steps.
- NP** All decision problems that can be solved by a *non-deterministic* turing machine using only polynomially-bounded many steps.
- L** All decision problems that can be solved by a *deterministic* turing machine using only log-bounded space for the computation.
- NL** All decision problems that can be solved by a *non-deterministic* turing machine using only log-bounded space for the computation.

Interestingly, similar to the open question $P = NP$ it is also still undecided whether $L = NL$.

To get a better understanding what you can do with a TM that has a logarithmic space bound, consider the following examples:

FIXME

Writing down all occurrences of a symbol in the input Intuitively one would argue that this problem requires a TM that visits each character

2.4 $NL \subseteq P$

In this section we want to assert where NL is placed in the complexity hierarchy. It should be clear that $L \subseteq NL$ and $P \subseteq NP$, so we are mainly interested in placing NL in that hierarchy. As it turns out, we can show that $NL \subseteq P$. Actually we can show much more:

Theorem 1. *Every turing machine that only needs logarithmic-bounded space for its computation and halts, also has a poly-bounded running time.*

Proof. A configuration of a TM is defined as the tuple (T, o, p, q) where $T : \mathbb{N} \rightarrow \Gamma$ is the current working tape state, $o, p \in \mathbb{N}$ which are the current head position on the input tape and on the working tape respectively and $q \in Q$ is the current state.

For a TM that only needs logarithmic-bounded space, we know that there are only $|\Gamma|^{O(\log n)} \in n^{O(1)}$ possibilities for T , as the number of symbols that can be non-blank is bounded by $O(\log n)$. Since the number of states is finite and the head position on the tape can be bounded by $O(\log n)$ as well (for the input tape it is always bounded by n), we can give a bound on the number of configurations:

$$n^{O(n)} \cdot n \cdot O(\log n) \cdot |Q| \subset O(n^k) \text{ for a suitably chosen } k \in \mathbb{N}.$$

Now we need to see that this results in a poly bounded running time for a NTM. Since we have an upper bound to the number of configurations, any cycle-free sequence of configuration that leads to a halting state is shorter than our bound. So a NTM only needs to iterate all configurations until a halting configuration is found, or the length of the sequence is longer than our bound. This can be done in polynomial running time thus, $NL \subseteq P$.

□

This proof is given for deterministic TM, since every non-deterministic TM can be transformed to an equivalent deterministic TM (with exponentially more states, but note that is still constant with regard to the input size).

2.5 Reducibility and NL-completeness

In the case of $P = NP$ it has been proven valuable to search for certain “hard” problems that characterize the complexity class. Similar to the notion of *NP-Completeness*, that is certain problems that are at least as difficult as any other decision problem in that complexity class, we try to define *NL-Completeness*.

A common technique to show that a decision problem is NP-complete, is to use reductions, for example the *Polynomial-Time Many-One Reduction*. A decision problem A is said to be *poly-time many-one reducible* to a decision problem B if and only if there exists a function f that can be computed in poly-time, for which the following requirement holds:

$$w \in A \Leftrightarrow f(w) \in B$$

Naively applying the same technique here will not work, since the poly-time constraint on the function is much too loose. Since the transforming function has no log-space constraint, it could be used to solve all decision problems in NL thus yielding a trivial reduction.

So, for NL-Completeness we need to add a space constraint on the transformation function: A decision problem A is said to be *log reducible* to B ($A \leq_{log} B$) if and only if there exists a function f that can be computed using logarithmic-bounded space (and thus is also poly-time constraint as we saw in the previous section) and satisfies the following requirement:

$$w \in A \Leftrightarrow f(w) \in B$$

We say B is *NL-complete* iff $A \leq_{log} B$ for all decision problems $A \in NL$.

Chapter 3

PATH and UPATH

The complexity class NL has a prominent member: the reachability problem in graphs. It actually turns out that this problem is different for directed and undirected graphs, which we will explore in this section.

3.1 PATH

Given a *directed* graph $G = (V, E)$ and two nodes $a, b \in V$ the decision problem *PATH* can be formulated as:

$(G, a, b) \in \text{PATH} \Leftrightarrow$ there exists a path from a to b in G .

3.1.1 NL-completeness

Theorem 2. *PATH is NL-complete.*

Based on [6] Solution to Exercise 5.3.

Proof. It is clear that $\text{PATH} \in \text{NL}$: Non-deterministically guess a path from a to b . This can be done by non-deterministically generating a sequence of n nodes in the graph and checking to see if adjacent pairs in the sequence are adjacent in the graph and if both a and b occur in the sequence.

Note: This works since each cycle-free path in a graph has at most n nodes, so if there is one the NTM will find it.

The only space needed on the working tape is the space to store a pair of nodes and a counter, which is all bounded by $\log n$.

Now let $A \in \text{NL}$ via a $O(\log n)$ space-bounded machine M . As we saw in the proof of Theorem 1, this machine has at most polynomially many configurations on an input of length n . The desired reduction of A to *PATH* outputs for any x the graph in which each such configuration is a node, and there is an edge from configuration c_i to c_j if c_j is a configuration that could follow c_i in the computation on input x . This can be done, for example, by producing for each configuration the finite list of possible successor configurations.

Note: This can be computed by a TM that also has logarithmic space-bound, even if the output is linear in the input since, since we do not consider the output tape!

By adding a single end-node, for all accepting configurations, we can use *PATH* to decide if there exists a path from the start configuration to the end-node. If so, we accept, otherwise we reject. \square

3.2 *UPATH*

Similar to *PATH* we define *UPATH* as:

$(G, a, b) \in \textit{PATH} \Leftrightarrow$ there exists a path from a to b in G .
where G is an *undirected* graph.

3.3 *UPATH* vs. *PATH*

One might be tempted to adapt the previous proof of Theorem 2 for NP-completeness of *PATH* to *UPATH*, but this will not work, since the given reduction will create invalid results for undirected graphs: It is often not possible to transition from c_i to c_j , even if you can transition from c_j to c_i . Thus *UPATH* would accept a lot of invalid computation paths in the undirected graph.

This might already be a hint, that *UPATH* is an easier problem than *PATH*, in fact we will see, that using randomization, we can solve *UPATH* using no non-determinism at all.

Chapter 4

RL

To differentiate *PATH* and *UPATH* even more, we will introduce yet another complexity class called *RL*.

RL contains all the decision problems that can be decided by a TM that only needs logarithmic-bounded space for the computation, but is allowed to execute a *random* decision at each step *and only uses polynomial-bounded many steps*.

Since we now have a *randomized* acceptor *M*, we need to redefine the notion of accepting:

$$\begin{aligned}x \in A &\Rightarrow \Pr[M \text{ accepts } x] \geq \frac{1}{2} \\x \notin A &\Rightarrow \Pr[M \text{ accepts } x] = 0\end{aligned}$$

Also, instead of running time and space usage, we use the *expected running time* and the *expected space usage*.

4.1 Randomized vs. Non-Deterministic

One might wonder what the relation between *RL* and *NL* is. It is important not to confuse the concepts of *random* (as used by *RL*) and *non-deterministic* (as used by *NL*). Non-determinism allows the TM to guess the *correct* transition in each step. A randomized TM can use a random value to determine the next transition, but for the *same random value* it will always act *deterministically*. Thus we see that every deterministic TM is a randomized TM that takes exactly zero random decisions. Furthermore every randomized TM can be simulated by non-deterministic TM, by replacing random decisions with non-deterministic transitions to the corresponding states that could be chosen based on the random value.

To conclude we see that: $L \subseteq RL \subseteq NL$

4.2 Proof that poly-running time bound is required

As you might have noticed, we required a poly-bound for running time of a TM in our definition of *RL*, which is something we could omit in the case of deterministic and non-deterministic TM. In this case however the requirement is non-optional. To understand why, we construct a randomized TM that has a log-space bound but has no polynomially bounded runtime.

Theorem 3. *There are randomized TM that have logarithmic space use, but an exponential running-time.*

Proof. Let *T* be a randomized turing machine, that 'flips a coin' with probability $p = \left(\frac{1}{2}\right)^n$. The first time it sees head it halts.

It should be clear that this TM does not need to use any space on the working tape.

To answer the question of the expected running time of such a TM, we model it as random experiment. X denotes the number of coin flips before the first head. The probability that we saw k tails before the first head is given by:

$$P(X = k) = (1 - p)^k \cdot p$$

which follows the geometric probability density. That yields:

$$E(X) = \frac{1}{p} - 1 = 2^n - 1$$

So the expected runtime is in $O(2^n)$

□

Chapter 5

Random Walk

Using the previous definition for TM that can use randomization we can execute the following algorithm, which we will call *RandomWalk*.

5.1 The algorithm

Input: (G, a, b)

```

 $v \leftarrow a$ 
for  $i \leftarrow 1$  to  $p(n)$  do
  Randomly select a node  $w$  that is adjacent to  $v$ 
   $v \leftarrow w$ 
  if  $v = b$  then
    accept.
  end if
end for
reject.

```

It should be clear, that this can be executed by a randomized TM M . But to show that $UPATH \in RL$ we still need to prove the following requirements:

1. M can compute *RandomWalk* using only logarithmically bounded space.
2. We can find a polynomial $p(n)$, for which M satisfies:

$$(G, a, b) \in UPATH \Rightarrow \Pr[M \text{ accepts } (G, a, b)] \geq \frac{1}{2}$$

$$(G, a, b) \notin UPATH \Rightarrow \Pr[M \text{ accepts } (G, a, b)] = 0$$

The first requirement is easy to prove, since we can see from the algorithm, that the only state that needs to be stored on the working tape is v . Since we can encode v using its node number, we need only $O(\log n)$ bits for that.

Proving the second requirement however, will compromise the most of the remainder of this chapter.

5.2 Application on PATH

At first sight, the above algorithm could also work for directed graphs, if one would define what happens in case *RandomWalk* reaches a node that has no outgoing edge. (e.g. simply restart at a) However we can show that there are directed graphs, for which doing only poly-bounded many steps will lead to incorrect results.

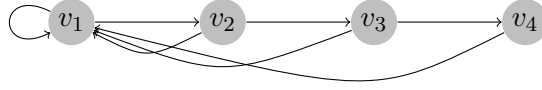


Figure 5.1: The graph used in proof of Theorem 4 for $n = 4$

Theorem 4. *There are directed graphs for which RandomWalk does not decide $(G, a, b) \in \text{PATH}$ correctly (as defined for randomized TMs) using only polynomially-bounded many steps.*

Proof. Since we want to be independent of defining what happens when *RandomWalk* reaches a dead-end, we simply provide a graph, where every node has an outgoing edge.

We construct a graph where for each node v_i there is an edge (v_i, v_{i+1}) and (v_i, v_1) for $n \in \{1, \dots, n\}$. We can see that $|E| = 2n$ and $|V| = n$, so the size of our input graph is in $O(n)$. We set $a = v_1$ and $b = v_n$.

See figure 5.1 for an example with $n = 4$.

At each node *RandomWalk* selects with $p = \frac{1}{2}$ the correct node. So the probability for finding the path from a to b is $p = (\frac{1}{2})^n$. Since each new try starting at a is independent from the previous one, this can be modelled as simple coin-flipping experiment using p as propability for success. As we saw in the proof of theorem 3, the expected number of retries will be 2^n .

Since the loop of the *RandomWalk* algorithm will do at least one iteration until we reach a again, the expected length of the random walk is at least 2^n . Thus *RandomWalk* can not decide instances of this class correctly using a polynomially-bounded many steps. \square

5.3 Correctness for UPATH

To show that our decider for UPATH based on *RandomWalk* really only needs to do poly-bounded many steps, we are going to abstract from the concrete algorithm and define a *random walk on a graph*.

A *random walk on a graph G starting at a* is a *infinite* sequence of nodes $W = (v_1, v_2, \dots)$ where each v_{i+1} was chosen randomly under uniform distribution from the neighbour nodes of v_i . (so $\{v_i, v_{i+1}\}$ is an edge in G)

We are now interested in the probability that a node v occurs in this sequence.

To do this, we are going to model the *random walk on G* as a *Markov Chain*. We define the random variable X_i as the node in G that is reached at the i -th step of our random walk. We can see that, X_i only depends on the value of X_{i-1} , the previous node, so we know that:

$$Pr[X_i = v | X_{i-1} = u] = \begin{cases} \frac{1}{d(u)} & \text{if } v \text{ is adjacent to } u \\ 0 & \text{otherwise} \end{cases}$$

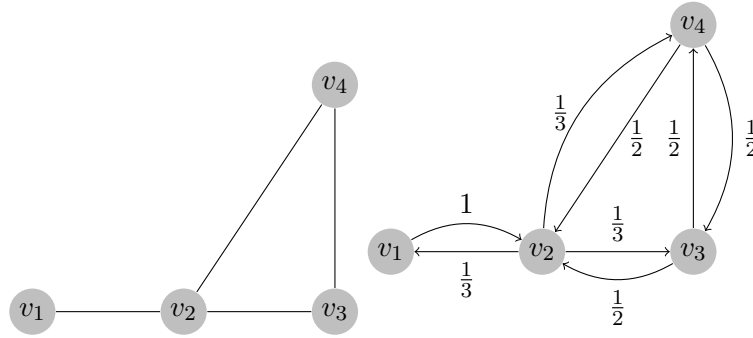


Figure 5.2: The original graph on the left, the state graph on the right.

Thus the sequence of random variables X_1, X_2, X_3, \dots has the *Markov Property*.

Furthermore we see that $\Pr[X_i = v | X_{i-1} = u]$ only depends on u and v and not on i from which we can conclude that this is a *stationary* markov chain.

This means it is rather easy to transform a graph into the state diagram of a markov chain that models a random walk on this graph: We replace the undirected edges $\{u, v\}$ by directed edges (u, v) and (v, u) and label them with the correct transition probabilities, $\frac{1}{d(u)}$ and $\frac{1}{d(v)}$ respectively. For an example see figure 5.2.

Since each input graph is of finite size, we can define the state-transition matrix $P \in \mathbb{R}^{n \times n}$

1. Show that $P_v = \lim_{n \rightarrow \infty} \frac{|\{i \leq n | v = v_i\}|}{n}$
 - Proof that $P_{(u,v)} = \frac{1}{2 \cdot e}$
 - Proof that $P_v = d(v) \cdot P_{(u,v)}$
2. Show that the expected number of steps of a random walk of a to visit all nodes once is bounded by a polynomial.
3. Show that b appears with probability higher than $\frac{1}{2}$ on such a random walk.

Chapter 6

Universal Traversal Sequence

6.1 Definition

- Visit all nodes in a d -regular Graph (why d -regular: formalism, node numbering..)

6.2 Relation to previous proof

- usage of probability application (only finitely many d -regular graphs of a certain size)

References

- [1] David J Aldous. Random walk covering of some special trees. *Journal of mathematical analysis and applications*, 157(1):271–283, 1991.
- [2] Romas Aleliunas, Richard M Karp, Richard J Lipton, Laszlo Lovasz, and Charles Rackoff. Random walks, universal traversal sequences, and the complexity of maze problems. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 218–223. IEEE, 1979.
- [3] Roy Armoni, Amnon Ta-Shma, Avi Wigderson, and Shiyu Zhou. An $o(\log(n)^{4/3})$ space algorithm for (s, t) connectivity in undirected graphs. *Journal of the ACM (JACM)*, 47(2):294–311, 2000.
- [4] Uriel Feige. A tight lower bound on the cover time for random walks on graphs. *Random Structures & Algorithms*, 6(4):433–438, 1995.
- [5] Michal Koucký. Universal traversal sequences with backtracking. *Journal of Computer and System Sciences*, 65(4):717–726, 2002.
- [6] Uwe Schöning and Randall Pruim. *Gems of theoretical computer science*. Springer, 1998.