

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	6
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
5.3.1	Modeling as Markov Chain	12
5.3.2	Computing a bound on the length of a random walk	15
5.3.3	Correctness of the decider	16
5.3.4	Conclusion	16
6	Universal Traversal Sequence	17
6.1	Definition	17
6.2	Finding the universal traversal sequence	17
	References	18

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 neighbor 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. In fact 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 minimum 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 separate 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 (Note: each cell can be visited multiple times, but will only be counted once). 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 example:

Writing down all occurrences of a symbol in the input Intuitively

one would argue that this problem requires a TM that needs at least $O(n)$ cells for the input and (in the worst case) $O(n \cdot \log n)$ cells for writing down the positions. Thus this problem has a space complexity of $\Omega(n)$. But since we only count cells that are visited on our *working tape*, the actual space requirement is the space we need for noting how often the symbol has occurred so far, which is in $O(\log n)$.

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 a non-deterministic TM, since every non-deterministic TM can be transformed to an equivalent deterministic TM (with exponentially more states, but note that it 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. Similarly 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 lax. 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 an input x the graph in which each 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:

$$x \in A \Rightarrow \Pr[M \text{ accepts } x] \geq \frac{1}{2}$$

$$x \notin A \Rightarrow \Pr[M \text{ accepts } x] = 0$$

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 running time.

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 $p = (\frac{1}{2})^n$. 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 occupy 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.

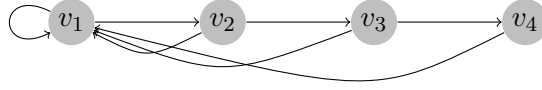


Figure 5.1: The graph used in proof of Theorem 4 for $n = 4$

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 modeled 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 an *infinite* sequence of nodes $W = (v_1, v_2, \dots)$ where each v_{i+1} was chosen randomly under uniform distribution from the neighbor nodes of v_i . (so $\{v_i, v_{i+1}\}$ is an edge in G)

We are now interested in the probability P_v that a node v occurs in this sequence. It should be clear that $P_v = 0$, if v_1 and v are not part of the same connected component. It is handy for the following chapter to assume that we have a connected graph.

5.3.1 Modeling as Markov Chain

To compute the value for P_v , 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

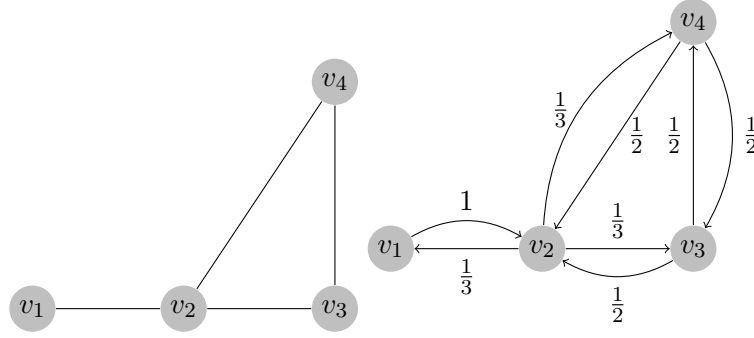


Figure 5.2: The original graph on the left, the state graph on the right.

random walk. We can see that, X_i only depends on the value of X_{i-1} , the previous node, so we know that:

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

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 (e.g. how many steps we have taken so far), from which we can conclude that this is a *time-homogenous* markov chain.

This means it is rather easy to transform a graph into the state-graph 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 $P \in \mathbb{R}^{n \times n}$ that computes for a distribution v the distribution after taking a step in our random walk. So P^n will compute the probability distribution resulting from taking n steps.

We now want to show the *unique* existence of a so called *stationary distribution*, that is a distribution π for which:

$$\pi \cdot P = \pi$$

Which means, we reached a distribution where taking further steps on our random walk does not change the probability that we are in a given state.

First, let us define what it means for a markov chain to be *irreducible*:

irreducibility If we take enough steps, eventually every transition (i, j) has a probability higher than zero.

Or more precise: $\exists i \in \mathbb{N} : Pr[X_i = v | X_0 = u] > 0$ for all $u, v \in V$

This simply follows from the observation that we can derive the state-graph from the *undirected connected* input graph.

Thus for each pair of states (u, v) there is a path from u to v , which will have probability > 0 .

Theorem 5. *For an irreducible markov chain with finite states, there is always a unique stationary distribution.*

Proof. The proof requires the definition of a lot of concepts, so we are only going to specify the general idea here. A complete proof can be found in any book about markov chains.

From the fact that we have finite many states, we see that for a markov chain of infinite size, we are going to visit at least one state infinitely often. This property is called *recurrence*.

If we know a recurrent state a , we can compute the expected value of how often we visit all other states before we return to a . In this case we can write that as vector ν . We can proof that $\nu \cdot P = \nu$.

If we norm ν we get our stationary distribution.

To see that it is unique, we show that there is one *positive recurrent* state ($\pi(i) > 0$ for at least one state i) and combined with the *irreducibility* we get that all states are *positive recurrent*.

A common theorem then specifies that each irreducible and positive recurrent markov chain has a unique stationary distribution. \square

Now that we know that there is a unique stationary distribution, we can proceed to find the actual values:

Theorem 6. *The stationary distribution for random walk is given by: $\pi = \frac{d(i)}{2e}$, where $e = |E|$.*

Proof. We simply show that the given value for π fulfills the identity $\pi P = \pi$.

It should be clear that $\sum_{i=0}^n \pi(i) = 1$ since $\sum_{i=0}^n d(i) = 2e$. (Remember: e is the number of *undirected* edges.) Furthermore:

$$\begin{aligned} \pi(i) &= \pi P_i \\ &= \sum_{j=1}^n \pi(j) \cdot \mathbf{1}_{\{i,j\} \in E} \cdot \frac{1}{d(j)} \end{aligned}$$

where P_i denotes the i -th column of P . Substituting $\pi(j) = \frac{d(j)}{2e}$ we get:

$$\begin{aligned} \pi(i) &= \sum_{j=1}^n \frac{d(j)}{2e} \cdot \mathbf{1}_{\{i,j\} \in E} \cdot \frac{1}{d(j)} \\ &= \frac{1}{2e} \cdot \sum_{j=1}^n \mathbf{1}_{\{i,j\} \in E} \\ &= \frac{1}{2e} \cdot d(i) \end{aligned}$$

\square

We can now use this stationary distribution to actually answer the question of a concrete value of P_v .

First we need to consider that the value of P_v is highly depended on how we choose the *start node*. We model choosing the start node as supplying an *initial distribution* to our markov chain.

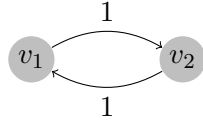


Figure 5.3: A simply cyclic markov chain. For a definition of the term cyclic we refer to appropriate literature.

Consider figure 5.3, which has $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$. If our *initial distribution* is $\nu = (1, 0)$ (that is we start at node v_1) we get $\nu \cdot P = (0, 1)$ and $\nu \cdot P^2 = (1, 0)$. So after two steps of our random walk we are back at v_1 with probability 1. It is clear we never reach our stationary distribution $(\frac{1}{2}, \frac{1}{2})$.

However if we assume we start at each node uniformly distributed with probability $\frac{1}{n}$ we see that for our cyclic example that would mean $\nu = (\frac{1}{2}, \frac{1}{2})$, which already is our stationary distribution.

It can be shown that this initial distribution will always lead to a convergence on our stationary distribution (follows from the Ergodic Theorem). To get an idea why this works, we can think of starting with an uniform initial distribution as conducting all possible random walks starting at all nodes at the same time. The probabilities we compute are the probabilities to be in state i over *all* random walks.

Now, finally we can see that $P_v = \pi(v)$ if we assume we start uniformly distributed. The expected number of steps between occurrences of a state in a markov chain is $\frac{1}{\pi(v)} = \frac{2e}{d(v)} =: E(v, v)$.

5.3.2 Computing a bound on the length of a random walk

To expand on that we are going to proof an upper bound for $E(i, j)$ with $\{i, j\} \in E$.

Theorem 7. $E(u, v) \leq 2e$ for $\{u, v\} \in E$.

Based on [6] Solution to Exerise 5.9.

Proof. The probability of chosing v if we are at node u is $\frac{1}{d(u)}$. That means the expected number of times we have to select one of u neighbours is $d(u)$, so we see:

$E(u, v) \leq E(u, u) \cdot d(u) = \frac{2e}{d(u)} \cdot d(u) = 2e$ Note that this is not a tight upper bound since v could be reached using different edges than (u, v) much earlier. \square

Using this we can proof an upper bound for $E(i, G)$, that is the expected number of steps of a random walk starting at v to visit *all* nodes. We can use this as upper bound for $E(i, j)$ (of course this bound is not exactly tight, but suffices).

First we should not that in an undirected connected graph, there is always a path (v_0, v_1, \dots, v_k) of length less than $2n$ that reaches all nodes. We are not going to state the proof for that here, but you can find it in [6] Solution to Exercise 5.10.

From that we can see that:

$$E(a, b) \leq E(a, G) \leq \sum_{i=1}^k E(v_{i-1}, v_i) < 2n \cdot 2e = 4en$$

5.3.3 Correctness of the decider

We are now going to proof that using $8en$ as an upper bound for the steps of our random walk starting at a we are *not* going to visit b with probability less than $\frac{1}{2}$ thus our decider will decide correctly with probability $\geq \frac{1}{2}$, which again means that our decider is a valid randomized decider.

Let $T(a, b)$ denote the random variable that models the number of steps a random walk from a to b takes. The probability that our random walk starting at a does not reach b using $8en$ steps, equals the probability that we need to take more than $8en$ steps to get from a to b .

Theorem 8. *Our random walk starting at a will not reach b with probability smaller than $\frac{1}{2}$:*

$$Pr[T(a, b) \geq 8en] < \frac{1}{2}$$

Proof. As we saw before, we have an upper bound for the expected value of $T(a, b)$ ($E(a, b) \leq 4en$). Since $T(a, b)$ is a positive random variable, we can apply the Markov inequality:

$$Pr[T(a, b) \geq 8e] \leq \frac{E(a, b)}{8e} < \frac{4en}{8en} = \frac{1}{2} \quad \square$$

5.3.4 Conclusion

As we saw in the beginning of the chapter, if a and b are not part of the same connected component, $P_b = 0$ which means we will never reach b for a random walk starting at a .

Using the theory of markov chains we can prove that we only need poly-bounded many steps $8en$ to find a path from a to b with probability higher than $\frac{1}{2}$.

Actually we saw (thanks to our rather generous upper bound) that using $8en$ steps, the probability that we *do not* reach all nodes in G is less than $\frac{1}{2}$. This insight is important for the next chapter on *Universal Traversal Sequences*.

Chapter 6

Universal Traversal Sequence

In this chapter we will step back from the decision problems *UPATH* and *PATH* and will discuss an interesting “by product” of our correctness proof in the previous chapter.

We saw that using a random walk with $8en$ steps we will traverse all nodes in a graph with probability higher than $\frac{1}{2}$. If we carry out m random walks we can decrease the probability that we do not reach all nodes: Let R_i denote the event that we do not reach all nodes using a random walk on G , then we see that:

$$Pr[R_1, \dots, R_m] = Pr[R_1] \cdot \dots \cdot Pr[R_m] < 2^{-m}$$

So our error decreases *exponentially*. However the number of steps $m \cdot 8en$ is still polynomial. We call this effect *probability amplification*.

We can use this effect to prove that there must be a sequence of directions that carries out a random walk on all d -regular graphs of the same size and visits all nodes (with probability 1.0).

6.1 Definition

First let us define what a *d-regular graph* is: We call a graph $G = (E, V)$ *d-regular* iff G is connected and $\forall v \in V : d(v) = d$. That is, each node has the same degree d .

The constraint on the degree of each node, makes the following definition possible:

We define a *traversal sequence* on such a graph as a sequence $I = (i_1, \dots, i_k)$ where each $i_q \in \{0, \dots, d-1\}$. If each node numbers its adjacent edges from 0 to $d-1$ this sequence is essentially a routing instruction for that graph. Note: This sequences are valid on *all* d -regular graphs, but most likely do not describe the same route.

6.2 Finding the universal traversal sequence

It is easy to see that in such a graph $|E| = e \leq \frac{d \cdot n}{2}$.

Since there is only a finite amount of d -regular graphs of a certain size, let us denote $g_{n,d}$ as the number of all d -regular graphs with n nodes.

Theorem 9. $g_{n,d} \leq n^{d \cdot n}$.

Based on [6] Solution to Exercise 5.12.

Proof. To build the adjacent edges of a node we need to chose d egdes. For each edge you can chose from at most n nodes, so that yields *at most* $n \cdot d$ possibilities to build the adjacent edges of a node. Since our graph has n nodes that results in *at most* n^{nd} possibilities to construct the edges for a d -regular graph with n nodes.

Thus $g_{n,d} \leq n^{d \cdot n}$. \square

If we choose i_q with uniform distribution for a traversal sequence $I = (i_1, \dots, i_k)$, we see that this sequence carries out a random walk of length k on each d -regular graph. As we saw before a random walk with $k = m \cdot 8en \leq 4dn$ has a probability of less than 2^{-m} to *not* visit all nodes.

Thus the expected number of d -regular graphs that a given traversal sequence will not traverse completely is at most $g_{n,d} \cdot 2^{-m}$. If we make m sufficiently large, such that $g_{n,d} \cdot 2^{-m} < 1$ there is a positive probability that a traversal sequence will traverse all nodes in all graphs. We call this sequence *Universal Traversal Sequence*.

Theorem 10. *There always exists a Universal Traversal Sequence and it has polynomially bounded length.*

Proof. First we need to find a m such that: $g_{n,d} \cdot 2^{-m} < 1$

$$\begin{aligned} g_{n,d} \cdot 2^{-m} &< 1 \\ n^{nd} \cdot 2^{-m} &< 1 \\ \log_2 n \cdot nd + (-m) &< 0 \\ \log_2 n \cdot nd &< m \end{aligned}$$

So there exists a traversal sequence with $k = 4dn \cdot (\log n \cdot nd + 1)$ that visits all nodes in all d -regular graphs of size n .

We see that $k \in O(n^3 \log n)$ (since $d \leq n$) thus the length is indeed bounded by a polynomial. \square

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 $o(\log(n) \cdot 4/3)$ 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.