

Course on Pseudorandomness

(Lecture Notes)

JAYALAL SARMA

Department of Computer Science and Engineering
Indian Institute of Technology Madras (IITM)
Chennai, India

Last updated on : February 17, 2021

Table of Contents

I	General Introduction and Tools	2
1	Week 1 : Power of Randomization and Derandomization Problem	3
1.1	Randomness helps in Matrix Multiplication Verification	3
1.2	Polynomial Identity Testing Problem	5
1.3	Derandomization Problem	8
1.3.1	Abstract Model of Derandomization	8
1.3.2	Derandomization by Brute Force Approach	8
1.4	Pseudorandomness : An Informal Overview	9
2	Week 2 : Method of Conditional Expectation, Pessimistic Estimators	13
2.1	MAXCUT Problem and Randomized Approximation	14
2.1.1	A Simple Randomized Algorithm	15
2.1.2	Recap of Probability Basics, Random Variables, Expectation	15
2.1.3	Analysis of the Algorithm for MAXCUT	21
2.2	Method of Conditional Expectation	21
2.2.1	Framework for Algorithms for Decision Problems	22
2.2.2	Framework for Algorithms for Optimization Problems	23
2.3	Method of Pessimistic Estimators	25
2.3.1	Congestion Minimization Problem	25
2.3.2	Randomized Approximation Algorithm	26
2.3.3	Pessimistic Estimator	27
3	Week 3 : Amplification : Making Algorithms Err Less	28
3.1	Success Probability Amplification by Repetition	29
3.2	Sipser's Argument	32
3.3	Amplification with Dependent Trials using Expanders	33
II	Exercise & Problem Sets	35
4	Exercises	36
4.1	Exercises	36
4.2	Curiosity Drive	37
5	Problem Sets	38
5.1	Problem Set #1	38

Todo list

1: Jayalal says: Todo - A few more lines to be completed here about the precise application of the estimators.	28
--	----

Part I

General Introduction and Tools

Randomized algorithms play a very powerful role in algorithm design. We will concentrate on the randomized algorithms for decision problems in this course. So all of our computational problems can be abstractly represented as given a string $x \in \Sigma^*$ in an alphabet, does x have property \mathcal{P} or not?

Informally, a randomized algorithm running in time t is an algorithm that on input x is allowed to perform at most t instances of random experiment of tossing unbiased coins during its computation and uses the outcome of the experiment in the computation, but however, provides a guarantee that the answer of the algorithm is the *correct* answer for the input x in a good fraction of the possible outcomes of the experiment. ¹

1.1 Randomness helps in Matrix Multiplication Verification

Consider the task of multiplying two $n \times n$ matrices over the field \mathbb{F}_2 . The trivial algorithmic solution to the problem takes $O(n^3)$ time and the trivial lower bound for the problem is $\Omega(n^2)$. It has been a long standing question which is the right complexity bound for this important problem (improvement to which will lead to improvements even in practice !). The exponent of matrix multiplication is the smallest constant ω such that two $n \times n$ matrices may be multiplied by performing $O(n^{\omega+\epsilon})$ for every $\epsilon > 0$.

Indeed, one of the basic ideas that we learn in algorithms courses for demonstrating the power of divide and conquer is the Strassen's multiplication which gives a running time bound of $O(n^{2.73})$ which went through a sequence of improvements and to the current best of $O(n^{2.31})$.

The question that we address now is something closely related - that of verifying whether a given multiplication is correct.

PROBLEM 1.1.1. Given three matrices $A, B, C \in \mathbb{F}_2^{n \times n}$, check whether $AB = C$ or not.

Indeed, the trivial method would be to multiply the two matrices and check if the result is equal to C . But then this requires, $O(n^{2.31} + n^2)$ time using the best matrix multiplication algorithm that we know currently. Since we just require verification, it is conceivable that we might be able to do better if we are allowed to make a small some error in the process. We show that this indeed possible to be done in $O(n^2)$ with error probability at most 2^{-k} for any constant k (independent of

¹Notice that if the guarantee for the algorithm is not saying "strictly more than half fraction of the coin tosses", then essentially the algorithm is useless since we can always replace it with a random experiment of tossing an unbiased coin and returning the answer to be YES if we get the heads and NO if we get tails. Note that we have at least a $\frac{1}{2}$ probability of success $\frac{1}{2}$.

n).

Trivial Approach using randomization: A natural first cut attempt is to choose an entry $(i, j) \in [n] \times [n]$ uniformly at random from the n^2 entries of C and checking if :

$$\sum_{k=1}^n A_{ik}B_{kj} = C_{ij}$$

This runs in time $O(n)$. And we can choose constant k more entries to amplify the success probability. However, the probability of correctness is very small. Suppose, in the worst case input, there was only one $(i, j) \in n \times n$ where there was an error in the multiplication. The probability that we will choose that particular (i, j) for verification is as small as $\frac{1}{n^2}$. Amplifying this to a success probability of $\frac{1}{2^k}$ takes more than $\Omega(n^{1+\epsilon})$ iterations and hence the overall algorithm will take $\Omega(n^{2+\epsilon})$ time which is beyond what we can afford to spend time on.

Freivalds' Approach: The idea is to check a randomly chosen "linear combination" of entries rather than a single entry of the matrix C . If we choose this to be a random linear combination of rows of the matrix C , then the combinatorics helps to achieve a much better probability of error. We now formally write down the algorithm and analyse it.

Algorithm 1.1 : Frievald's Algorithm for Verification of Matrix Multiplication - $\mathcal{F}(A, B, C)$

- 1: Choose a vector $r \in \mathbb{F}_2^n$ uniformly at random.
 - 2: If $[A(Br) = Cr]$ then output YES else output NO.
-

The computation of $(A(Br))$ and Cr are done using $O(n^2)$ time algorithms since computing a linear transformation result Ax for an $n \times n$ matrix can be done in $O(n^2)$ time. Now we argue correctness guarantees. If the given matrices indeed satisfy $AB = C$, then no matter which $r \in \mathbb{F}_2^n$ algorithm chooses in step 1, $ABr = Cr$ and hence it always will output YES. The error can happen only when $AB \neq C$ and the algorithm ends up choosing an unfortunate r such that $ABr = Cr$. The following claim upper bounds the probability of this .

CLAIM 1.1.2. For any $A, B, C \in \mathbb{F}_2^n$ such that $AB \neq C$,

$$\Pr[\mathcal{F}(A, B, C) \text{ outputs YES}] \leq \frac{1}{2}$$

Proof. Let $A, B, C \in \mathbb{F}_2^n$ such that $AB \neq C$. We need to analyze the probability that $ABr = Cr$ for $r \in \mathbb{F}_2^n$ chosen uniformly at random. If $AB \neq C$, then $D = AB - C$ is a non-zero matrix. Thus

$$\Pr_{r \in \mathbb{F}_2^n} [ABr \neq Cr] = \Pr_{r \in \mathbb{F}_2^n} [Dr \neq 0]$$

Imagine that D was all 1s matrix. Now the above probability is exactly the number of vectors $r \in \mathbb{F}_2^n$ with an odd number of 1s in it. By an obvious bijection, this is also the number of subsets of $[n]$ with odd size. The latter is exactly 2^{n-1} and hence the probability of such a vector r being chosen from \mathbb{F}_2^n is $\frac{1}{2}$.

Now we formalize and generalize this. Let p be the vector Dr . Since $D \neq 0$, there must be an entry $D_{ij} \neq 0$. Define, $A = \{j : D_{ij} \neq 0\}$. We know that $A \neq \emptyset$. $i, j \in [n]$. Thus :

$$p_i = \sum_{k=1}^n D_{ik} r_k = \sum_{k \in A} r_k$$

$$\text{Note that: } \Pr_{r \in \mathbb{F}_2^n} [Dr = 0] \leq \Pr_{r \in \mathbb{F}_2^n} [p_i = 0]$$

Notice that the latter probability depends only on r_k where $k \in A$. The fraction of assignments of the bits $\{r_k : k \in A\}$ which makes $p_i = 0$ is exactly $\frac{1}{2}$ since the number of even sized subsets of A and number of odd sized subsets of A are exactly the same. \square

REMARK 1.1.3. Informally, if we run the above algorithm $\mathcal{F}(A, B, C)$, and the algorithm outputs NO, then we can trust the answer and conclude that indeed $AB \neq C$. But a YES answer from the algorithm cannot be trusted - it could be because of the unfortunate choice of $r \in \mathbb{F}_2^n$ that came as the outcome of the experiment.

Why are we interested in the above algorithm even though it gives a success probability bound of only $\frac{1}{2}$? The reason is that, it is a one-sided error algorithm and hence still much better than a coin toss outcome because the algorithm does not make an error when $AB = C$. In fact, such algorithms can be repeated in a natural way - run k times, and if any of them says $AB \neq C$ output NO. This reduces the error probability in exponentially in k - since each of the trials should give an error (with probability $\frac{1}{2}$) and hence the error probability bound is at most $\frac{1}{2^k}$.

1.2 Polynomial Identity Testing Problem

The previous example, while it demonstrates the point, might be a bit unsatisfactory since the problem under consideration anyway has an efficient algorithm. To address this, we will now see another example problem where there is an efficient (polynomial time in the input size) randomized algorithm for solving the problem but a deterministic algorithm for solving the problem is not known.

The problem is easy-to-state algorithmic question on polynomials. Fix \mathbb{F} to be the field where the coefficients are chosen from. *Given a polynomial $p \in \mathbb{F}[x_1, x_2, \dots, x_n]$, test if it is identically zero.* That is, do all the terms cancel out and become the zero polynomial.

This problem has its roots in the simple high school arithmetic. Suppose we are given a polynomial in a complicated form where the monomials may repeat with arbitrary coefficients etc. We want to find out if the coefficient of the monomials cancel out to zero. This in effect is testing whether the polynomial is the zero polynomial, and equivalently it is testing if the polynomials evaluates to zero on all substitutions of the variable from the underlying field \mathbb{F} .

How are we given the polynomial? This indeed is going to have effect on the complexity of the problem. Let us start with the high school arithmetic again. Suppose we are given it in the monomial form (though some monomials may repeat) along with their coefficients. To solve the problem, it suffices to check, for each monomial whether the coefficient in its various appearances is adding up to zero. Given the explicit representation at the input, this is very easy to do by simply going over the input for each monomial. Hence this can be done in time polynomial in the

input.

What if the polynomial is not given that explicitly. How can it be given implicitly compared to the list of monomials? One answer is that, we could give it in a bracketed form. That is, the polynomial $x_1x_2 + x_1x_4 + x_3x_2 + x_3x_4$ can be given as $(x_1 + x_3)(x_1 + x_4)$. Indeed, this is implicit, since an expression of length n , can have number of actual number of monomials to be 2^n - consider the example $(x_1 + x_2)(x_3 + x_4) \dots (x_{n-1} + x_n)$ where n is the number of variables.

What is the most implicit form that we can think of? A black box which evaluates the polynomial. That is, we have an oracle p when given input a returns $p(a)$, the value of polynomial at a .

Assume that we are also given an upper bound on the degree of the polynomial $\deg(p) \leq d$. Indeed, we do not have access to the actual polynomial except through the blackbox. We have to use some property of the degree d polynomials. The most obvious one is the number of points in which they can evaluate to zero. Based on this thought, the following deterministic algorithm solves the problem.

Algorithm 1.2 A deterministic algorithm for univariate polynomial identity testing

- 1: Choose $d + 1$ different points a_1, \dots, a_{d+1} .
 - 2: Call the oracle $d + 1$ times to evaluate $p(a_1), \dots, p(a_{d+1})$.
 - 3: If all calls returned 0 accept else reject.
-

If p were really the zero polynomial then all calls will return 0 and we will definitely accept. If p were not 0, then at most d calls can return 0 since a polynomial with degree at most d has at most d roots. Hence if $p \neq 0$, then our algorithm will definitely reject.

Multivariate PIT: Now let us think about the problem when p is a multivariate polynomial. The previous assertion that a degree d polynomial has at most d roots no longer holds. To see this, consider the degree 2 polynomial $p(x_1, x_2) = x_1x_2$. This has an infinite number of roots $x_1 = 0, x_2 \in \mathbb{F}$, where \mathbb{F} is the (possibly infinite) field over which p is defined.

We can work around this problem by considering a finite subset of the field, say $S = \{0, \dots, 10\}$. The polynomial p has 19 zeroes. So if x_1, x_2 is chosen uniformly at random from S there is at most 19/100 chance that we will get a false result. As can be seen from the above example, by making the size of S arbitrarily large, we can make the error probability arbitrarily small. But then the disadvantage is that we will need more random bits in order to choose an element at random from the set $|S|$, and the running time of our algorithm will also increase.

Generalizing this strategy that we will follow is as follows: If the total degree of the polynomial is $\leq d$, and if $S \subseteq \mathbb{F}$, such that $|S| \geq 2d$, instead of picking elements arbitrarily, we pick elements uniformly at random from S . Indeed, there may be many choices for the values which may lead to zero. But how many?

LEMMA 1.2.1 (Schwartz-Zippel Lemma). *Let $p(x_1, x_2, \dots, x_n)$ be a non-zero polynomial over a field \mathbb{F} . Let $S \subseteq \mathbb{F}$*

$$\Pr_{\vec{a} \in S^n} [p(\vec{a}) = 0] \leq \frac{d}{|S|}$$

Proof. (By induction on n) For $n = 1$: For a univariate polynomial p of degree d , there are $\leq d$

roots. Now in the worst case the set S that we picked has all d roots. Thus for a random choice of substitution for the variable from S , the probability that it is a zero of the polynomial p is at most $\frac{d}{|S|}$.

For $n > 1$, write the polynomial p as a univariate polynomial in x_1 with coefficients as polynomials in the variables $p(x_2, \dots, x_n)$.

$$\sum_{j=0}^d x_1^j p_j(x_2, x_3, \dots, x_n)$$

For example: $x_1 x_2^2 + x_1^2 x_2 x_3 + x_3^2 = (x_2 x_3) x_1^2 + (x_2^2) x_1 + x_3^2$.

We need to analyze the probability that we will choose a zero of the polynomial (even though the polynomial is not identically zero). For a choice of the variables as $(a_1, a_2, \dots, a_n) \in S^n$, we ask the question : how can $p(a_1, a_2, \dots, a_n)$ be zero? It could be because of two reasons:

1. $\forall j : 1 \leq j \leq n, p_j(a_2, a_3, \dots, a_n) = 0$.
2. Some coefficients $p_j(a_2, a_3, \dots, a_n) = 0$ are non-zero, but the resulting univariate polynomial in x_1 evaluates to zero upon substituting $x_1 = a_1$.

Now we are ready to calculate $\Pr[p(a_1, a_2, \dots, a_n) = 0]$. For a random choice of (a_1, \dots, a_n) . Let A denote the event that the polynomial $p(a_1, \dots, a_n) = 0$. Let B denote the event that $\forall j : 1 \leq j \leq n, p_j(a_2, a_3, \dots, a_n) = 0$. Note that, $\Pr[A] = \Pr[A \wedge B] + \Pr[A \wedge \bar{B}]$.

We calculate both the terms separately: $\Pr[A \wedge B] = \Pr[B].\Pr[A|B] = \Pr[B]$ where the last equality is because $B \Rightarrow A$. Let ℓ be the highest power of x_1 in $p(x)$. That is $p_\ell \neq 0$. Since the event B insists that for all j , $p_j(a_2, a_3, \dots, a_n) = 0$, we have that $\Pr[B] \leq \Pr[p_\ell(a_2, a_3, \dots, a_n) \neq 0]$. By induction hypothesis, since this polynomial has only $n - 1$ variables and has degree at most $\frac{d-\ell}{5}$. Thus, $\Pr[B] \leq \frac{d-\ell}{5}$.

To calculate the other term,

$$\Pr[A \cap \bar{B}] = \Pr[\bar{B}].\Pr[A|\bar{B}] \leq \Pr[A|\bar{B}] \leq \frac{\ell}{|S|}$$

where the last inequality holds because the degree of the non-zero univariate polynomial after substituting for a_2, \dots, a_n is at most ℓ and hence the base case applies. \square

This suggests the following efficient algorithm for solving PIT. Given d and a blackbox evaluating the polynomial p of degree at most d .

Algorithm 1.3 : Schwartz-Zippel Algorithm for Multivariate PIT

- 1: Choose $S \subseteq \mathbb{F}$ of size $\geq 4d$.
 - 2: Choose $(a_1, a_2, \dots, a_n) \in_R S^n$.
 - 3: Evaluate $p(a_1, a_2, \dots, a_n)$ by querying the blackbox.
 - 4: If it evaluates to 0 accept else reject.
-

The algorithm runs in time $\text{poly}(n)$. The following Lemma states the error probability and follows from the Schwartz-Zippel Lemma that we saw before.

LEMMA 1.2.2. *There is a randomized polynomial time algorithm A , which, given a black box access to a polynomial p of degree d (d is also given in unary), answers whether the polynomial is identically zero or not, correctly with probability at least $\frac{3}{4}$.*

Notice that in fact the lemma is weak in the sense that it ignores the fact that when the polynomial is identically zero then the success probability of the algorithm is actually 1 !. In other words, is it is a one-sided error randomized algorithm.

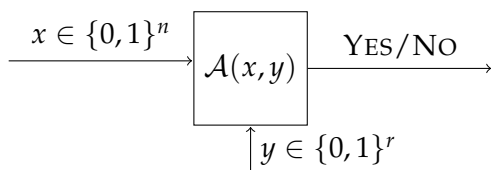
1.3 Derandomization Problem

In the previous section (lecture), we talked about randomized algorithms for problems for which we do not know deterministic algorithms with similar complexity resource bounds. Indeed, we are not happy about randomized algorithms as such since these algorithms require perfect unbiased coin toss experiments to be performed and we do not have them in practice. Indeed, the fact that they can output erroneous answers, even though with low probability makes them useless in critical practical applications.

How do convert them to deterministic algorithms without causing much overhead?. One possible way is to look at each algorithm and use inherent properties of the problem to analyze the randomized algorithm better to come up with ways to remove randomness from that algorithm. Here, we start with the original randomized algorithm for a particular problem, and improve it to derandomize it and the techniques are usually very algorithm specific. We will do some examples of this kind later in the course.

1.3.1 Abstract Model of Derandomization

From now on, we will be concentrating only on abstract models of these randomized algorithms. We fix some notations first. A randomized algorithm \mathcal{A} on input x runs in time $t(n)$ (where $n = |x|$) and let $y \in \{0, 1\}^{r(n)}$ be the concatenation of the unbiased coin toss experiment that the algorithm does during its execution. Notice that $r(n) \leq t(n)$ (we drop the n when it is not required explicitly). If the algorithm runs in polynomial time $t(n) \leq n^c$ for a constant c independent of n .



The guarantee we have is there is an $\epsilon \in (0, \frac{1}{2}]$.

$$\forall x \in \{0, 1\}^n, \Pr_{y \in \{0, 1\}^r} [A(x, y) \text{ is correct.}] \geq \frac{1}{2} + \epsilon$$

1.3.2 Derandomization by Brute Force Approach

The trivial approach to obtain an equivalent deterministic algorithm is run over all possible outcomes of the experiment and check the answer from the algorithm for each of them. Whichever answer comes as majority - report that as the final answer.

Algorithm 1.4 (\mathcal{A}') : input $x \in \{0,1\}^n$, where success prob. $\frac{1}{2} + \epsilon$ for \mathcal{A}

```
1:  $count \leftarrow 0$ .
2: for each  $y \in \{0,1\}^r$  do
3:   Check if  $\mathcal{A}(x,y)$  accepts, if so increment  $count$ 
4: end for
5: If  $[count > 2^{r-1}]$  then output YES else output NO.
```

If the running time of the randomized algorithm \mathcal{A} is $t(n)$, then the running time of the new algorithm (which is deterministic) is $t(n)2^{r(n)}$. To argue correctness, if the actual answer for input $x \in \{0,1\}^n$ is YES, then the fraction of $y \in \{0,1\}^r$ which makes $\mathcal{A}(x,y)$ accept is strictly more than $\frac{1}{2}$ and hence the algorithm will output YES. If the actual answer for input $x \in \{0,1\}^n$ is NO, then the fraction of $y \in \{0,1\}^r$ which makes $\mathcal{A}(x,y)$ accept is strictly less than $\frac{1}{2}$ and hence the algorithm will output NO.

REMARK 1.3.1. Note that the algorithm \mathcal{A} will run in $\text{poly}(n)$ time if the original randomized algorithm was running in $t \leq \text{poly}(n)$ time and was using $r \leq O(\log n)$ random bits.

1.4 Pseudorandomness : An Informal Overview

Ideally, we would like to replace the randomized algorithm with a deterministic one as done in the previous section. However, we know how to do this trivially only when the randomized algorithm uses $O(\log n)$ random bits.

We outline two “out of the box” thoughts related to our target of derandomization of randomized algorithms.

Fooling the Algorithm with Pseudorandom bits : PRGs The first one is about using $y \in \{0,1\}^r$ as not independent random bits. But use *dependent* random bits instead. Indeed, the analysis for the error bound for the algorithm \mathcal{A} now may fail since it may assume total independence between the bits of y in its mathematical argument. However, sometimes, it is possible that same analysis (or even a better analysis) may work even when the bits of the y are dependent in a limited way ² But this may be specific to the algorithm and sometimes to the problem itself. We would ideally want a more abstract strategy which would work for randomized algorithms in general, modelled by what we described in the previous section. But even if we made it work with some dependent randombits, how do we produce this distribution of y' with the desired limited dependence among them? Construction of the methods which can produced limited dependence thus becomes important.

Taking a more abstract view point, informally, we would like to have a box (formally an algorithm G) which takes in pure random bit string of length $y' \in \{0,1\}^{r'}$ and produces a string $y \in \{0,1\}^r$ such that the distribution of y “looks” pseudo-random for the resource limited algorithm \mathcal{A} . The idea is that we will run the algorithm \mathcal{A} with the random string provided the $G(y')$

²At one extreme, if we had an algorithm and an analysis which works wen all the bits of the y are the same (which is an example of extreme dependence) we dont require the random bit at all - we can directly simulate the algorithm deterministically the trivial way.

- the output of G on input $y' \in \{0,1\}^{r'}$ chosen uniformly at random.

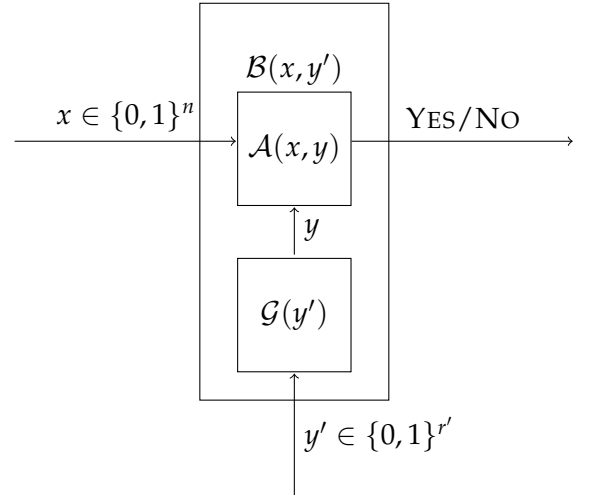
Indeed, since we are not providing pure random bits to \mathcal{A} . Hence we should expect its correctness guarantees to not hold good anymore. That is, it will deteriorate a bit. Can we guarantee that it does not deteriorate too much? This can be done in two ways (1) rework or reanalyse the algorithm \mathcal{A} and argue that it is still having success probability greater than $\frac{1}{2}$ (which is enough for trivial derandomization) (2) use only resource bounds of \mathcal{A} to argue that the change of y to $G(y')$ will not affect the success probability much. For this, the generator has to satisfy certain properties.

A function $G : \{0,1\}^{r'} \rightarrow \{0,1\}^r$ is said to be a Pseudo-Random Generator (PRG) for complexity measure³ s and error parameter $\delta \in [0,1]$ if, for any algorithm \mathcal{A} which runs in time $t \leq s$ (or having complexity measure bounded by s): For any x ,

$$\left| \Pr_{y \in \{0,1\}^r} [\mathcal{A}(x,y) \text{ Accepts}] - \Pr_{y' \in \{0,1\}^{r'}} [\mathcal{A}(x, G(y')) \text{ Accepts}] \right| \leq \delta$$

Connecting to the informal description, δ is the quantity by which the success probability deteriorates because of the use of the pseudorandom generator output, instead of pure random bits. Hence if we ensure that $\delta < \epsilon$, even after the use of the pseudorandom generator output, we still will have a randomized algorithm \mathcal{B} with the following guarantee $\forall x \in \{0,1\}^n$:

$$\Pr_{y \in \{0,1\}^r} [\mathcal{B}(x,y) \text{ is correct.}] \geq \frac{1}{2} + \epsilon - \delta > \frac{1}{2}$$



We will end the description by asking the question - *what parameters determine how good our pseudorandom generator is?*. As per the above discussion it is:

- The relative values of r and r' . This leads to the definition of the *stretch* of the pseudorandom generator. We would ideally want an exponential stretch function so that with $r' \in O(\log n)$ we can produce y for \mathcal{A} which is of length $r = O(n^c)$ for constant c .
- The value of s . This determines how powerful an algorithm can the pseudorandom generator manage to fool. The larger the s the better. Ideally we want s to be covering all polynomial time running time bounds.
- The value of δ . This determines the quantity by which the success probability of the algorithm deteriorates after plugging in the output of the pseudorandom generator instead of the y from the pure random bits. Ideally, we want $\epsilon - \delta > 0$. The smaller the δ , the better.

³We will make it more precise when it comes to the section where we handles these objects. We are leaving at the above description at a less precise level.

- Running time of the generator itself. Notice that we need \mathcal{G} to be explicit polynomial time algorithm, which runs in time $\text{poly}(n)$. That is, if $r' \in O(\log n)$ (which is what ideally we would want, so that the trivial derandomization runs in $\text{poly}(n)$ time), then technically, the generator can run for exponential time in terms of its input size⁴.

The main part of the game is in describing the generator algorithms (or functions from $\{0,1\}^{r'} \rightarrow \{0,1\}^r$). However, it is not even clear whether such functions exists for the range of parameters that we care about. Indeed, this is the kind of flavour that we will have.

- We can prove that the functions that we are looking for exist, with a non-constructive argument. This is done by - what is termed as the *Probablistic method*.
- Explicit descriptions of the functions, which are are required for the algorithms with required runtime bounds for \mathcal{G} are not known. In fact, if we have such descriptions, then a complete derandomization of all randomized algorithms is possible, which will be a big achievement.

Refining Randomness : Randomness Extractors - Here is a completely different idea about supplying dependent randomness. We do not have source of pure random bits to supply for the randomized algorithm \mathcal{A} . But we may have impure random bit sources. An imaginative question is *can we invest a few pure random bits in order to purify/extract and hence improve the impurity in the given random source?*. This, at first sounds crazy and leads to the following questions.

- How do we define *impure random bit sources*. They define distributions which are not uniform. There is the notion of entropy which can tell us how uniform the source is.
- How do we define *how good the output is*. Again, one could have used entropy here too naturally. However, noticing the fact that we would like to finally apply it our algorithms like \mathcal{A} , a different measure of "purity" is used which is the notion of statistical distance⁵ to uniform distribution.

The above discussion leads to the definition of a randomness extractor, which is a function $\mathcal{E} : \{0,1\}^n \times \{0,1\}^d \rightarrow \{0,1\}^m$ such that when X is a distribution on $\{0,1\}^n$ with entropy at least k , then the distribution of the output $\mathcal{E}(X, U_d)$ is ϵ -close to U_m where U_d and U_m are uniform distributions on the set $\{0,1\}^d$ and $\{0,1\}^m$ respectively.

Again, how do we determine how good is our extractor function? We want extractors which works on highly biased distributions (the smallest k possible) using fewest number of pure random bits (the smallest d possible) and produces output distributions which are closest to uniform distributions (ϵ must be smallest) - and still run in time polynomial in n .

Similar to pseudorandom generator functions, it is unclear apriori whether such functions even exist for the range or parameters we care about. A similar situation arises, where by using probabilistic method, we can prove that such objects (functions) exists, but at the same time, we do not know how to construct them deterministically (equivalently describe the algorithm for \mathcal{E}).

⁴This marks the difference between the pseudorandom generators studied in cryptography and derandomization.

⁵Informally, this is the sum of the difference (in absolute value) between the probability values assigned to points in the sample space.

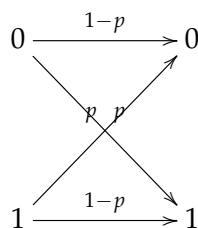
REMARK 1.4.1 ((Informal) - **Psuedo-random Objects**). The common flavour that we observed about the previous sections is that there are mathematical objects which we would like to describe (by providing algorithms to compute those functions) and we do know that the function that we seek exists. This situation is a common phenomenon in many objects. In fact, in most situations, it is not just the existence of the objects that is argued, but also that if the object that is of interest is chosen at random (with appropriately set up experiments), the object of interest shows up at the outcome with high probability. Thus, there is a randomized algorithm to explicitly construct the object, and now we have to derandomize them !. However, notice that in such situations, we can think of derandomizations which just depends on those algorithms which chooses the object at random.

Other Contexts and Psuedorandom Objects We now describe a totally unrelated context in which the required mathematical functions display such a psuedorandom behaviour and an explicit construction is being sought for. The context is that of coding theory.

Coding theory had its inception in the late 1940's with the theory of reliable communication over a channel in the presence of noise - an area that started with the pioneering work of Claude Shannon and Richard Hamming. The former addressed and answered the fundamental questions about the possibility of the use of codes for reliable communication and the later developed some basic combinatorial constructions of error correcting codes that laid the foundations for the work later.

Theoretical computer scientists have a major role to play in the algorithmic aspects of coding theory research, and coding theory has proved to be instrumental in several interesting results in theoretical computer science as well. There has been several surprising applications of codes and the associated mathematical objects, in areas like algorithms, complexity theory and cryptography. Some part of the course will aim to discuss of those applications. However, we do not intend to be exhaustive.

The channel is not harmless in the real world. It introduces errors in the transmission. Depending on the application the error may be in the physical storage media (communication over time) or in the physical channel (communication over space). Some of the 0s gets flipped to 1s and vice versa, and some bits may get dropped too. For the purposes of this course we will study only model (Shannon studied several interesting variants), namely what are called Binary Symmetric Channels. In this model, each bit gets flipped with a probability p . That is, a 1 gets flipped to a 0 with probability p and 0 gets flipped to 1 with probability p .



What is the natural strategy to cope up with errors in transmission? Create redundancy. For example, if Alice wants to send a bit 0 to Bob, she will do it five times, and send 11111 and ask Bob to take the majority of the bits as the bit that was sent. In this simple looking example we

have all the essence. The string that was sent will be called the *codeword* and the original bit to be sent is called the *message*. There are only two codewords 00000 and 11111 in the above example. If we define the notion of distance as the hamming distance, then the majority decoding mechanism described above can also be seen as choosing the codeword that is closest to the received word. This natural strategy of decoding is called *nearest neighbor decoding* or *maximum likelihood decoding*.

Now let us observe facts about guarantees. Clearly if the channel is such that it will not corrupt more than 2 bits in a sequence of 5 bits, then Bob will be able to decode the message bit correctly. But the channel may actually flip more number of bits but with relatively lower probability. Thus if we increase the number of copies we make of the original message, with high probability (over the errors) introduced by the channel we are going to be able to decode the bit correctly.

To fix some notations, we denote $E : \{0,1\}^k \rightarrow \{0,1\}^n$ as the encoding function where k is the message length (in general) and n is the length of the codeword (which we will call the *block length*). Let $m \in \{0,1\}^k$ be a message, and $E(m) \in \{0,1\}^n$ is the transmitted word. The channel corrupts the message and let $y \in \{0,1\}^n$ is the received word. The error introduced by the channel could also be thought of as a string $\eta \in \{0,1\}^n$ where the η_i determines whether $y_i = (E(m))_i$ or not.

We want the following guarantee for any $m \in \{0,1\}^k$ as translating the above intuition:

$$\Pr_{\eta}(D(E(m) + \eta) = m) \geq 1 - o(1)$$

where the $o(1)$ term is exponentially small depending on n and hence on k (since c is a constant).

Although the above statement is written in terms of a probability over choice of the channel error vector, a natural combinatorial guarantee that we would want is an encoding and decoding scheme such that if the error string η has weight at most $t < \frac{d}{2}$ the decoder retrieves the message correctly. That is, the encoder-decoder pair is guaranteed to get the message across the channel, if the number of corruptions by the channel is limited a number t . Indeed, the relative redundant information we sent should be minimised (which is the ratio of k and n called the rate of the code).

Shannons theorem essentially states that under suitable choice of the parameters there is a pair of encoding-decoding functions that can achieve this high confidence decoding of the original message. We will state the theorem formally only later. But again, the spirit of the theorem is that there does exist good encoding and decoding schemes with respect to the parameters we usually care about (which we make precise later). The area of algorithmic coding theory essentially attempts to address the question of constructing coding schemes for which there is an efficient decoding.

We conclude the lecture by stating that the three mathematical objects that we stated in this lecture do have some interconnections among themselves and also the pseudorandom objects that we are going to state in the next lecture too.

In this week, the plan is to learn the technique of conditional expectation which is a derandomization technique that works for several algorithms. However, this is a technique which is specific to algorithms and not to problems. And there is no hard and fast rule by which we can say whether this technique is applicable for an algorithm. We demonstrate it with the MAXCUT problem.

2.1 MAXCUT Problem and Randomized Approximation

For an undirected graph $G(V, E)$, a cut is a partition of vertices into two sets $S, T \subseteq V$. The size of a cut is the number of edges that go across these partitions. That is, the size of the set :

$$\text{cut}(S, T) = \{e = (u, v) \mid (u, v) \in E, u \in S, v \in T\}$$

Maximum cut is a cut whose size is at least the size of any other cut. That is $|\text{cut}(S, T)|$ is the largest possible. Given a graph, the problem of finding a maximum cut in a graph is known as the MAXCUT problem.

The problem is hard: The MAXCUT problem is known to be NP-hard. This implies, in particular, that if we have an efficient algorithm for the MAXCUT problem, then some of the very hard problems will yield to having efficient algorithms solving them. This is believed to be unlikely.

Approximation algorithms: Hence it makes sense to talk about algorithms which may not output the exact maximum cut, but instead another cut. Indeed, this is useless unless there are guaranteed how large is the cut output by the algorithm. An example guarantee that we may want to target is, for the algorithm, no matter what the input graph G is, the cut output by the algorithm will be, say, at least $(\frac{1}{10})$ -th of the size of the maximum cut. This is called a 0.1-approximation algorithm⁶. Even with this relaxed target for the algorithm, is not immediately clear how to design such an algorithm. It turns out that the best known algorithm for MAXCUT does much better than this and achieves an approximation guarantee of 0.875, and for many reasons this is believed to be the best possible ratio that any polynomial time algorithm can achieve for MAXCUT problem. However, in this lecture, we will concentrate on much smaller ratios.

Randomized Approximation algorithms: We resort to randomized algorithms - which in this context will be called randomized approximation algorithms. Notice that unlike the previous

⁶Exercise: if you have not seen it already, think about what would be a similar statement that you would like to target for a minimization problem, like vertex cover problem

examples, MAXCUT is not a decision problem. Hence, we need to be careful about designing and analysing randomized algorithms for it. For example, there is nothing like the algorithm being correct. Instead, we have only the notion of the approximation ratio - that is how close the output of the algorithm is, to the optimal value.

2.1.1 A Simple Randomized Algorithm

We start with a simple randomized algorithm for MAXCUT problem.

Algorithm 2.5 : Randomized Approx. Algorithm for MAXCUT for graph $G(V, E)$, $|V| = n$

```

1:  $S = T = \phi$ 
2: for each  $i \in [n]$  do
3:   Choose bit  $b_i \in \{0, 1\}$  uniformly at random.
4:   if  $b_i = 1$  then
5:      $S = S \cup \{i\}$ 
6:   else
7:      $T = T \cup \{i\}$ .
8:   end if
9: end for
10: Output  $cut(S, T)$ .
```

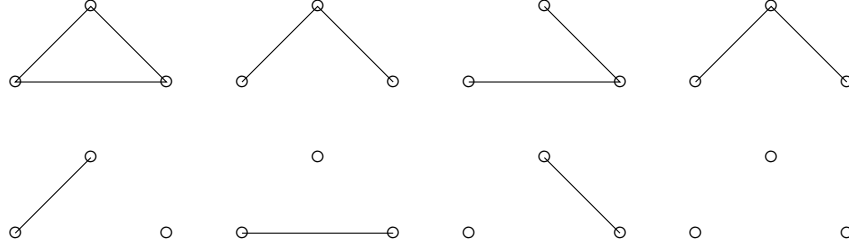
Notice that the sets S and T will form a partition of V at the end of the algorithm. Clearly, the algorithm will run in linear time. Indeed, the above description is also equivalent to - choose a random subset S of vertices from V and outputting $cut(S, \bar{S})$. Indeed, we need to give a guarantee about the size of the cut output by the algorithm. Roughly, the claim is that the average size of the cut (where average is taken over the 2^n different outcomes of the random choices). We recap some basics of random variables and expectations now before stating the correctness claim.

2.1.2 Recap of Probability Basics, Random Variables, Expectation

Fix a set Ω , which is called the *sample space*. A probability distribution is defined over Ω , is a function $\Pr : \Omega \rightarrow [0, 1]$ satisfying the additional condition that, $\sum_{w \in \Omega} \Pr(w) = 1$. An event is a subset $\mathcal{E} \subseteq \Omega$. The probability of an even \mathcal{E} is nothing buy the sum of the probability values by assigned by the distribution functon to the elements in the subset \mathcal{E} . That is, $\Pr(\mathcal{E}) = \sum_{w \in \mathcal{E}} \Pr(w)$.

We consider an example which are going to be relevant for us. This is the notion of random graphs. Consider n vertices, and there are $\binom{n}{2}$ possible edges. Imagine that we will choose (independently) each edge to be present in our graph with probability p and to be absent in our graph with probability. The outcome of the experiment is an n -vertex simple graph and hence the sample space is the set of all n vertex graphs.

As an example, we can consider, $n = 3$. There are 3 possible edges and hence 8 possible graphs. The probability assigned the triangle graph (which is the complete graph on 3 vertices) is p^3 since all the three edges have to be chosen for this particular outcome to happen. In a similar way, the following pictures denote the sample space in this case with the corresponding probability values.



The probabilities in that order are p^3 , $p^2(1-p)$, $p^2(1-p)$, $p^2(1-p)$, $p(1-p)^2$, $p(1-p)^2$, $p(1-p)^2$, $(1-p)^3$.

How do we analyse the probability that we get a connected graph as the outcome of the experiment. This is where events are used. Recall that formally, an event \mathcal{E} is a subset of Ω .

$$Pr(\mathcal{E}) = \sum_{w \in \mathcal{E}} Pr(w)$$

In the above example, if the event \mathcal{E} represent the set of connected graphs.

$$Pr(\mathcal{E}) = p^3 + 2p^2(1-p)$$

In the above example, if the event \mathcal{E}' represent the set of bipartite graphs.

$$Pr(\mathcal{E}') = 1 - p^3$$

PROPOSITION 2.1.1 (Subadditivity of Probability - a.k.a - Union theorem). Let $\mathcal{E}_1, \mathcal{E}_2 \dots \mathcal{E}_n$ be events, then :

$$Pr \left[\bigcup_i \mathcal{E}_i \right] \leq \sum_{i=1}^n Pr[\mathcal{E}_i]$$

DEFINITION 2.1.2 (Conditional Probability). For two events \mathcal{E} and \mathcal{E}' , we define,

$$Pr(\mathcal{E}|\mathcal{E}') = \frac{Pr(\mathcal{E} \cap \mathcal{E}')}{Pr(\mathcal{E}')}$$

The conditional probability captures the questions of the kind, what is the probability that we get a connected graph if we are given that the outcome is a bipartite graph?

DEFINITION 2.1.3 (Independent events). Two events \mathcal{E} and \mathcal{E}' are said to be independent, if

$$Pr(\mathcal{E}|\mathcal{E}') = Pr(\mathcal{E})$$

Equivalently,

$$Pr(\mathcal{E} \cap \mathcal{E}') = Pr(\mathcal{E})Pr(\mathcal{E}')$$

For example, if we consider the events event \mathcal{E} represent the set of connected graphs and event \mathcal{E}' represent the set of bipartite graphs, then:

$$Pr(\mathcal{E} \cap \mathcal{E}') = 3p^2(1-p)$$

$$Pr(\mathcal{E})Pr(\mathcal{E}') = [(1-p)^3 + 3p(1-p)^2 + 3p^2(1-p)](1-p^3)$$

Since they are not equal, we conclude that the two events are not independent. That is, the event that the graph is bipartite has an "influence" on the event that the graph is connected. To make this clearer, we suggest the following exercise:

Exercise 2.1.4. Let $G \in G(n, p)$. For all $S \subseteq V$, let A_S be the event that S forms an independent set in G . Show that if S and T are two distinct subsets of k vertices then A_S and A_T are independent if and only if $|S \cap T| \leq 1$.

Now, we will generalize the above notion of independence to more than two events. An event \mathcal{E} is independent of a set of events $\{\mathcal{E}_j \mid j \in J\}$ if, for all subset $J' \subseteq J$, $Pr[\mathcal{E} \mid \cap_{j \in J'} \mathcal{E}_j] = Pr(\mathcal{E})$.

Exercise 2.1.5. Prove that an event \mathcal{E} is independent of a set of events $\{\mathcal{E}_j \mid j \in J\}$ if and only if for all $J_1, J_2 \subseteq J$ such that $J_1 \cap J_2 = \emptyset$

$$Pr[\mathcal{E} \cap (\cap_{j \in J_1} \mathcal{E}_j) \cap (\cap_{j \in J_2} \overline{\mathcal{E}_j})] = Pr(\mathcal{E}) Pr[(\cap_{j \in J_1} \mathcal{E}_j) \cap (\cap_{j \in J_2} \overline{\mathcal{E}_j})]$$

Let $\{\mathcal{E}_i \mid i \in I\}$ be a (finite) set of events. They are *pairwise independent* if for all $i \neq j$ the events \mathcal{E}_i and \mathcal{E}_j are independent. Events are *mutually independent* if each of them is independent from the set of the others. It is important to note that events may be pairwise independent but not mutually independent. Following exercise demonstrates that.

Exercise 2.1.6 (See Problem Set 1(Problem 1)). A random k -colouring for a graph G is an element of the probability space (Ω, Pr) where Ω is the set of all k -colourings (i.e. partition of V into k sets (V_1, V_2, \dots, V_k) , all this colourings being equally likely (so happening with probability $\frac{1}{k^n}$). For every edge e of G , let A_e be the event that the two endvertices of e receive the same colour. Show that:

- (a) for any two edges e and f of G , the events A_e and A_f are independent.
- (b) if e, f and g are three edges of a triangle of G , the events A_e, A_f and A_g are dependent.

Random Variables: We need the idea of random variables which we recap now. A random variable is another function $X : \Omega \rightarrow \mathbb{R}$. The expected value of the random variable is the "weighted average" value that it takes over the real numbers - weighted by the corresponding probability values. That is,

$$E[X] = \sum_{\alpha \in \mathbb{R}} \alpha Pr[X = \alpha]$$

Indeed, $[X = \alpha]$ represents an event $\{w \in \Omega \mid X(w) = \alpha\} \subseteq \Omega$. Hence, the expectation can also be written equivalently as follows:

$$E[X] = \sum_{\alpha \in \mathbb{R}} \alpha \left(\sum_{\substack{w \in \Omega \\ X(w) = \alpha}} Pr(w) \right) = \sum_{w \in \Omega} X(w) Pr(w)$$

We need the following properties of expectation:

Tool 1 : Boolean Random Variables - Suppose X is a random variable that takes only Boolean values. In this case, $E[X] = \Pr[X = 1]$ which follows from the definitions.

Tool 2 : Linearity of Expectation : Suppose X_1 and X_2 are random variables defined based on the same probability distribution, consider the new random variable defined as $X = c_1 X_1 + c_2 X_2$. This is also a random variable as it is a function from $\Omega \rightarrow \mathbb{R}$ defined as $X(w) = c_1 X_1(w) + c_2 X_2(w)$ for every $w \in \Omega$. It turns out there is a neat relationship between the expectation of the random variables X, X_1 and X_2 . This is one of the most important relation that is extensively used in analysis of randomized algorithms.

$$\begin{aligned} E[X] &= \sum_{w \in \Omega} X(w) \Pr(w) = \sum_{w \in \Omega} (c_1 X_1(w) + c_2 X_2(w)) \Pr(w) \\ &= c_1 \left(\sum_{w \in \Omega} X_1(w) \Pr(w) \right) + c_2 \left(\sum_{w \in \Omega} X_2(w) \Pr(w) \right) = c_1 E[X_1] + c_2 E[X_2] \end{aligned}$$

We suggest practicing the application of linearity of expectation using the following exercise.

Exercise 2.1.7 (See Problem Set 1(Problem 2)). A graph $G = (V, E)$ is created at random by selecting each edge with probability p . What is the expected number of spanning trees in the randomly sampled graph? (Hint : Use Cayley's Theorem that the number of distinct spanning trees on n vertices is n^{n-2} . Order them, and define an indicator random variable.)

Tool 3 : Averaging Principle - Suppose X is a random variable and $E[X] = \mu$, then the following statements follow:

$$\exists w \in \Omega : X(w) \geq \mu \quad \exists w \in \Omega : X(w) \leq \mu$$

Both of them can be proved by contradiction. For sample, we spell out the first one, suppose that the first statement is false. That is, $\forall w \in \Omega, X(w) < \mu$, then:

$$E[X] = \sum_{w \in \Omega} X(w) \Pr(w) < \sum_{w \in \Omega} (\mu \times \Pr(w)) = \mu \left(\sum_{w \in \Omega} \Pr(w) \right) = \mu$$

This implies, $E[X] < \mu$ which is a contradiction. A similar proof holds for the other claim as well.

Tool 4 : Tail inequalities - Suppose we have a random variable X such that $E[X] = \mu$. What kind of probability guarantees can we write for X ? For example, can we bound (in terms of the expectation) the probability that $X > \alpha$ for some $\alpha \in \mathbb{R}$? This is what tail bounds do. They help us write probability upper bounds based on expectations and other related parameters. As a first example, consider a random variable that takes only non-negative values. Then we can write :

$$\textbf{Markov's Inequality} : \Pr[X \geq a] \leq \frac{E[X]}{a}$$

The proof is also quite simple.

$$E[X] = \sum_{\alpha \in \mathbb{R}} \alpha \Pr[X = \alpha] \geq \sum_{\alpha \geq a} \alpha \Pr[X = \alpha] \geq \sum_{\alpha \geq a} a \Pr[X = \alpha] \geq a \Pr[X \geq a]$$

For example, this helps us make statements of the form :

In particular, Markov's inequality implies the following bound on the probability that any random variable X is significantly larger than its expectation:

$$\Pr[X \geq (1 + \delta)E[X]] \leq \frac{1}{1 + \delta} \quad (2.1)$$

For example, the probability that the random variable takes a value which is more than 4 times the expected value is at most 0.25. Unfortunately, this does not help us write down a probability bound for X taking value less than say $\frac{E[X]}{4}$. Indeed, another form is :

Indeed, Markov's inequality is also pretty weak - as demonstrated by the following example - consider tossing n coins and X be the number of heads. Clearly $E[X] = \frac{n}{2}$. By Markov's inequality, $\Pr[X \geq n] \leq \frac{1}{2}$ but we know that it is much smaller than that, namely $\frac{1}{2^n}$.

What do we do if we want to bound the probability that the random variable takes a much lower value than the expectation? This is where we require more the next tail bound (without any assumption of positivity on the random variable).

Chebychev's Inequality : $\Pr[|X - \mu| \geq a] \leq \frac{\text{Var}[X]}{a^2}$ where, $\text{Var}[X] = E[X^2] - E[X]^2$

Proof of Chebychev's Inequality in the sum of random variables case: More often, we require the Chebychev's inequality when X is a random variable that is expressible as a sum of several independent random variables. We will handle that case. In particular, we will prove that:

$$\Pr[(X - \mu)^2 \geq a] \leq \frac{E[X^2]}{a} \quad (2.2)$$

$X = \sum_{i=1}^k X_i$. Let $E[X_i] = \mu_i$ and $E[X] = \sum_i \mu_i = \mu$ (say).

Define $Y_i = X_i - \mu_i$. By linearity of expectation, $E[Y_i] = 0$. And let $Y = \sum_i Y_i = X - \mu$. The idea is to apply Markov's inequality to a positive random variable that can be formed out of Y . Indeed Y can take negative values, hence we can consider Y^2 . We care about $E[Y^2]$ in order to apply the Markov's inequality.

$$E[Y^2] = E\left[\left(\sum_i Y_i\right)\left(\sum_i Y_i\right)\right] = E\left[\sum_{i,j} Y_i Y_j\right] = \sum_{i,j} E[Y_i Y_j] = \sum_i E[Y_i^2] + \sum_{i \neq j} E[Y_i Y_j]$$

Now we will apply independence of the random variables, and conclude that $E[Y_i Y_j] = E[Y_i]E[Y_j]$ but then in our case $E[Y_i] = 0$. Using this:

$$E[Y^2] = \sum_i E[Y_i^2] = \sum_i [\mu_i(1 - \mu_i)^2 + (1 - \mu_i)\mu_i^2] = \sum_i \mu_i(1 - \mu_i) < \mu$$

Hence, the theorem (equation 2.2) follows from Markov's inequality. The general form also has a similar proof. And note that $\text{Var}[Y] = \mathbb{E}[Y^2]$ since $\mathbb{E}[Y] = 0$. \square

Chebychev's inequality gives significantly better bounds compared to Markov's inequality. For example, it implies (compare with Equation 2.1).

$$\Pr[X \geq (1 + \delta)\mu] \leq \frac{1}{\delta^2\mu} \quad (2.3)$$

REMARK 2.1.8. Note that we did not use independence fully - we only needed that for all $i \neq j$, $\mathbb{E}[Y_i Y_j] = \mathbb{E}[Y_i]\mathbb{E}[Y_j]$. This is a strictly weaker condition than saying all random variables Y_1, Y_2, \dots, Y_k are independent. Such random variables are called "pairwise independent". We will come across them later in the course.

We now present one of the stronger tools to estimate probability tails in the case of sum of fully independent random variables. In this case, the set up itself is about a random variable $X = \sum_{i=1}^k X_i$. Let $\mathbb{E}[X_i] = \mu_i$ and $\mathbb{E}[X] = \sum_i \mu_i = \mu$ (say). We have the following:

Chernoff Bound : For any $\alpha > \mu$, $\Pr[X \geq \alpha] \leq e^{\alpha - \mu} \left(\frac{\mu}{\alpha}\right)^\alpha$

$$\text{In particular, it implies: } \Pr[X \geq (1 + \delta)\mu] \leq \left(\frac{e^\delta}{(1 + \delta)^{(1 + \delta)}}\right)^\mu$$

Proof of Chernoff Bound: Again the idea is to use Markov's inequality for an appropriately defined positive random variable. In this case, we will just use $Y = \left(\frac{\alpha}{\mu}\right)^X$ and it is related to the original probability that we wanted to estimate. More precisely,

$$\Pr[X \geq \alpha] = \Pr\left[\left(\frac{\alpha}{\mu}\right)^X \geq \left(\frac{\alpha}{\mu}\right)^\alpha\right] = \Pr\left[Y \geq \left(\frac{\alpha}{\mu}\right)^\alpha\right] \leq \frac{\mathbb{E}[Y]}{\left(\frac{\alpha}{\mu}\right)^\alpha}$$

Since $\mathbb{E}[Y]$ is $\mathbb{E}\left[\left(\frac{\alpha}{\mu}\right)^X\right]$, we have the motivation to estimate $\mathbb{E}[a^X]$ in terms of $\mathbb{E}[X]$. Here we use the fact that $X = \sum_i X_i$ and are fully independent. Note that $\mathbb{E}[a^X] = \prod_i \mathbb{E}[a^{X_i}]$. Hence, we need to estimate $\mathbb{E}[a^{X_i}]$.

$$\begin{aligned} \mathbb{E}[a^{X_i}] &= a^0 \Pr[X_i = 0] + a^1 \Pr[X_i = 1] \\ &= 1 \cdot (1 - \mathbb{E}[X_i]) + a \cdot \mathbb{E}[X_i] = a\mu_i + (1 - \mu_i) = \mu_i(a - 1) + 1 < e^{(a-1)\mu_i} \end{aligned}$$

In our context $a = \frac{\alpha}{\mu}$ and $Y = \left(\frac{\alpha}{\mu}\right)^X$. This gives,

$$\Pr[X \geq \alpha] \leq \frac{\mathbb{E}[Y]}{\left(\frac{\alpha}{\mu}\right)^\alpha} \leq \frac{\mathbb{E}\left[\left(\frac{\alpha}{\mu}\right)^X\right]}{\left(\frac{\alpha}{\mu}\right)^\alpha} \leq \frac{\prod_i \mathbb{E}\left[\left(\frac{\alpha}{\mu}\right)^{X_i}\right]}{\left(\frac{\alpha}{\mu}\right)^\alpha} \leq \frac{\prod_i e^{(a-1)\mu_i}}{\left(\frac{\alpha}{\mu}\right)^\alpha} \leq \frac{e^{(\frac{\alpha}{\mu}-1)\mu}}{\left(\frac{\alpha}{\mu}\right)^\alpha} \leq e^{(\alpha-\mu)} \left(\frac{\mu}{\alpha}\right)^\alpha$$

□

2.1.3 Analysis of the Algorithm for MAXCUT

Recall Algorithm ???. We want to guarantee that the expected size of the cut is at most half of the optimal cut size. In fact, we prove something stronger.

CLAIM 2.1.9. *Let X be the size of the cut output by the algorithm ???. Then, $E[X] \geq \frac{m}{2}$ where m is the number of edges in the graph G .*

Proof. For each edge $e \in E$ define a random variable X_e as the following indicator variable:

$$X_e = \begin{cases} 1 & \text{if } e \in \text{cut}(S, T) \\ 0 & \text{otherwise} \end{cases}$$

By definition, $X = \sum_{e \in E} X_e$. Hence, by linearity of expectation:

$$E[X] = \sum_{e \in E} E[X_e] = mE[X_e]$$

We just need to notice that $E[X_e] = \Pr[X_e = 1] = \Pr[e \in \text{cut}(S, T)]$ because of tool 1. Notice that an edge e is in the cut if the two end points get into different sets among S and T . That is, out of the four possible outcomes of the random coin tosses corresponding to the endpoint vertices of e , two of them leads to e being in $\text{cut}(S, T)$. Hence this is exactly $\frac{1}{2}$. This gives: $E[X] \geq \frac{m}{2}$. Hence the proof. □

Notice that the claim is stronger. Indeed, since the optimum cut can only cut at most m edges, the above also implies $E[X] \geq \frac{m}{2} \geq \frac{\text{OPT}_{\text{CUT}}}{2}$.

2.2 Method of Conditional Expectation

We now describe the main technical idea to be learned this week. Mainly an algorithm specific technique of derandomization of randomized algorithms. This presentation is from Salil Vadhan's book on Pseudorandomness.

There are two kinds of randomized algorithms that we have seen so far - essentially to solve two kinds of problems. One is for the decision problems where the probability over different paths of the computation tree, the algorithm being correct is at least $\frac{2}{3}$. The other is for optimization problems where the expected size of the output has guarantees. The method of derandomization that we are going to discuss can in principle be applied for both, if the randomized algorithm in question uses the random bits in a peculiar way that certain measures can be computed efficiently about the output for particular settings of the randomness.

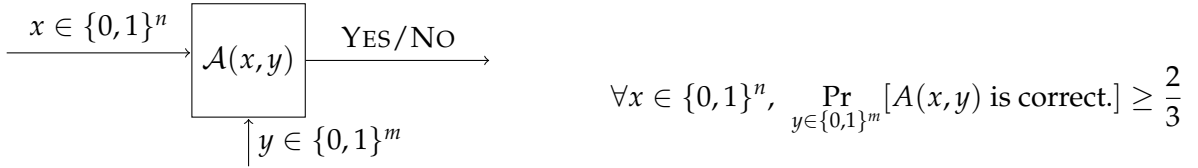
Main Idea: (as we described in the lecture) - we discuss the first kind of algorithms first and then adapt it to the second type. In the decision problem case, we know that $\frac{2}{3}$ -rd paths in the computation tree are going to make the algorithm answer correctly. A vague idea would be - walk down the path of the tree, *making a choice deterministically and efficiently at each node (without trying*

both choices which leads to exponential time) maintaining the invariant that within the subtree that have restricted ourselves to, a $\frac{2}{3}$ fraction of paths within the subtree still make the algorithm go correct. If we make choices like this and set the random bits based on that choices, it is intuitive that we will reach a leaf that makes the algorithm answer correctly and then we can just run the algorithm on that leaf (that choice of random bits) and output the answer. The process is deterministic and efficient and hence gives a derandomization of the original algorithm.

Of course, this is easier said than done. An important question remains - *at an intermediate node, in the above walk-down, how do we deterministically and efficiently decided whether to take the edge labelled 0 or 1 to move to the child?* Formalizing this will require us to fix the type of problem as decision vs search/optimization problem.

2.2.1 Framework for Algorithms for Decision Problems

Recall the following notational set up where \mathcal{A} is the randomized algorithm solving a decision problem.



Since we have to keep track of the fraction of paths for a particular partial setting of the random bits (while analysing the walk-down of the tree at an intermediate stage) - we define the following notation:

For every $i \in [n]$, bits $r_1, r_2, \dots, r_i \in \{0, 1\}$, define:

$$p(r_1, r_2, \dots, r_i) = \Pr_{y \in \{0, 1\}^m} \{ \mathcal{A}(x, y) \text{ is correct} \mid (y_1 = r_1) \wedge (y_2 = r_2) \wedge \dots \wedge (y_i = r_i) \}$$

Indeed, if we are at a particular node represented by a partial assignments r_1, r_2, \dots, r_i , the value of $p(r_1, r_2, \dots, r_i)$ is the average of the value at the two children of that node since the bit is chosen uniformly at random. In terms of expectation, this is equivalent to :

$$p(r_1, r_2, \dots, r_i) = E_{y_{i+1} \in \{0, 1\}} (r_1, r_2, \dots, r_i, y_{i+1})$$

To understand this definition clearly, let us ask, a first question, what is the value of $p(r_1, r_2, \dots, r_m)$ for a setting $r_1, r_2, \dots, r_m \in \{0, 1\}$. (Class answered 0 or 1 depending on whether the setting represents a path which makes \mathcal{A} correct or not). How about $p(\phi)$, which represents the value of the function when no bit is set. Clearly, by definition, this represents the top of the computation tree, and hence the fraction of correct paths under the node is exactly the success probability of the algorithm. That is, $p(\phi) \geq \frac{2}{3}$.

Indeed, if we call $p(r_1, r_2, \dots, r_i) = \mu$, we have that: $E_{y_{i+1} \in \{0, 1\}} p(r_1, r_2, \dots, r_i, y_{i+1}) = \mu$. Hence we know that there must exist a setting of y_{i+1} such that the value of the random variable - which

in this case is $p(r_1, r_2, \dots, r_i, y_{i+1})$ - is at least the expected value. That is,

$$\exists r_{i+1} \in \{0, 1\} : p(r_1, r_2, \dots, r_{i+1}) \geq p(r_1, r_2, \dots, r_i)$$

Applying this repeatedly, we have that $\exists r_1, r_2, \dots, r_m \in \{0, 1\}$:

$$p(r_1, r_2, \dots, r_m) \geq p(r_1, r_2, \dots, r_{m-1}) \geq \dots \geq p(r_1, r_2) \geq p(r_1) \geq p(\phi) \geq \frac{2}{3}$$

Notice that, by our observation, the left-end term is Boolean, and hence it must be that there exists $r_1, r_2, \dots, r_m \in \{0, 1\}$ such that $p(r_1, r_2, \dots, r_m) = 1$. But then, this is not a big deal in the end, we knew about existence of such r_i 's anyway. So in the end it does not look very useful.

But the above framework has an interesting feature. It also shows how to construct r_i 's bit-by-bit, if we have an efficient algorithm to compute $p(r_1, r_2, \dots, r_i)$ for any i . Suppose we have computed r_1, \dots, r_i already, we can compute r_{i+1} as follows: compute $p(r_1, r_2, \dots, r_i, 0)$ and $p(r_1, r_2, \dots, r_i, 1)$ using the above algorithm and set r_{i+1} to be whichever bit in $\{0, 1\}$ which achieves the maximum.

However, we still have the problem of computing $p(r_1, r_2, \dots, r_i)$ for any $i \in [m]$ efficiently. This is where the algorithm-specifics come in. The algorithm \mathcal{A} should be using the random bits in a peculiar way such that the value of $p(r_1, r_2, \dots, r_i)$ can be computed efficiently for that algorithm \mathcal{A} . Indeed, the trivial method of computing $p(r_1, r_2, \dots, r_i)$ ends up taking exponential time in the worst case. Hence one has to use the algorithm-specific attributes to design this algorithm. We will demonstrate this in the next context.

2.2.2 Framework for Algorithms for Optimization Problems

We now adapt the above framework for search and optimization problems. The guarantee for algorithms solving such problems is as follows - the expected size of the output is at least as "good" as this, where "good" means at most or at least in minimization and maximization problems respectively. For demonstrative purposes, we restrict ourselves to the randomized algorithm for MAXCUT that we presented earlier. Notice that the algorithm uses exactly n random bits where $|V| = n$. Note that S and T are the two subsets output by the algorithm. For any $i \in [n]$, define:

$$V(r_1, r_2, \dots, r_i) = \mathbb{E}_{y \in \{0, 1\}^n} [|cut(S, T)| : (y_1 = r_1) \wedge (y_2 = r_2) \wedge \dots \wedge (y_i = r_i)]$$

Similar to the previous setting, note that $V(\phi) \geq \frac{|E|}{2}$ since the expected size of the cut when no random bit is conditioned is similar to the original analysis of the algorithm. We apply the averaging principle and argue in a similar way that there must exist a choice of the random bits $r_1, r_2, \dots, r_n \in \{0, 1\}$ such that $cut(S, T)$ output by the algorithm has at least $\frac{|E|}{2}$ many edges.

Indeed, we start with $V(\phi) \geq \frac{|E|}{2}$. By averaging principle, there must exist $r_1 \in \{0, 1\}$ such that $V(r_1) \geq V(\phi) \geq \frac{|E|}{2}$. Continuing in a similar way there must exist $r_1, r_2, \dots, r_n \in \{0, 1\}$ such that :

$$V(r_1, r_2, \dots, r_{n-2}, r_{n-1}, r_n) \geq V(r_1, r_2, \dots, r_{n-2}, r_{n-1}) \geq V(r_1, r_2, \dots, r_{n-2}) \dots \geq V(r_1) \geq V(\phi) \geq \frac{|E|}{2}$$

Again, this is not new information. We can always derive by a globally applying averaging principle that there must exist such a choice of random bits. But the advantage here is that if there is an efficient algorithm for computing $V(r_1, r_2, \dots, r_i)$ for any i , then we can find out the explicit choice of values of the bits as well. As in the previous case, if r_1, r_2, \dots, r_i is already fixed, then we compute r_{i+1} as : compute $V(r_1, r_2, \dots, r_i, 0)$ and $V(r_1, r_2, \dots, r_i, 1)$ and set r_{i+1} to be that value in $\{0, 1\}$ which results in the maximum among the two.

Computing $V(r_1, r_2, \dots, r_i)$ for the algorithm for MAXCUT : To apply the above framework, all we need is an efficient algorithm to compute the value of $V(r_1, r_2, \dots, r_i)$ for any choice of i and $r_1, r_2, \dots, r_i \in \{0, 1\}$. Note that this is a speciality of the algorithm - more importantly the way the bits y_1, y_2, \dots, y_n are used by the algorithm.

At any intermediate point of computation, there are vertices for which the decision (of whether they should be in the set S or not) is already made by then and there are vertices which are decided later. To keep track of this, we define:

$$\begin{aligned} S_i &= \{j \in [n] \mid j \leq i, b_j = 1\} \\ T_i &= \{j \in [n] \mid j \leq i, b_j = 0\} \\ U_i &= \{j \in [n] \mid j > i\} \end{aligned}$$

The algorithm will grow S_i to S and T_i to T , by randomly choosing the remaining vertices (U_i) to be in S or T . We need to compute the expected size of the cut conditioned on the fact that the sets S_i and T_i are already fixed by the algorithm. An immediate observation is that the edges that go across S_i and T_i will necessarily a part of the cut, since their endpoints are already at S and T respectively by definition. The edges which are fully within S_i or fully within T_i are not going to be a part of the cut finally. But there may be more number of edges which forms a part of the final cut. Considering this, we can write:

$$V(r_1, r_2, \dots, r_i) = |cut(S_i, T_i)| + \frac{|(cut(S_i, U_i)| + |cut(U_i, T_i)| + |cut(U_i, U_i)|}{2}$$

We need to explain the second term in the RHS. Consider edges $e = (u, v) \in E$ that has one endpoint $u \in S_i$ and the other endpoint in $v \in U_i$. Note that $u \in S$ finally, and hence (u, v) edge will be counted in the cut, if v falls into T . Since this is decided by choosing a random bit, the probability that the edge appears in the cut finally is $\frac{1}{2}$. Hence, the expected number of edges in $cut(S_i, U_i)$ which appear in the final cut is $\frac{|cut(S_i, U_i)|}{2}$. Similar argument explains the term $\frac{|cut(T_i, U_i)|}{2}$. To see the $\frac{|cut(U_i, U_i)|}{2}$ term, consider edges which are having both end points in U_i . They are both going to be put in S or T uniformly at random - hence out of the four possible outcomes (for these two vertices), two of them puts them in the final cut and two of them puts them outside the final cut output by the algorithm. Hence for any edge in $cut(U_i, U_i)$, with probability $\frac{1}{2}$ it will form a part of the cut. That is, expected number of edges that gets contributed to the final cut is $\frac{|cut(U_i, U_i)|}{2}$. Hence the expression for $V(r_1, r_2, \dots, r_i)$ is correct.

Exercise 2.2.1 (See Problem Set 1(Problem 3)). In the derandomization of MAXCUT algorithm that we described, we derived an expression for $V(r_1, r_2, \dots, r_i)$ for any i and $r_1, r_2, \dots, r_i \in \{0, 1\}$. We

used this to determine, the value of r_{i+1} by computing $V(r_1, r_2, \dots, r_i, 0)$ and $V(r_1, r_2, \dots, r_i, 1)$ and then choosing the value of r_{i+1} to be the one which produces the largest among the two. Prove that the choice of r_{i+1} will be 1 if vertex $i + 1$ has more neighbors in T_i than in S_i and vice versa. Hence, write down the derandomized 0.5-approximation deterministic polynomial time algorithm for MAXCUT as a simple greedy algorithm in terms of the above rule.

Exercise 2.2.2 (See Problem Set 1(Problem 4)). Let $x_1, x_2, \dots, x_n \in \{0, 1\}$ be Boolean variables and let f be a Boolean formula in CNF form. That is, $f = C_1 \wedge C_2 \wedge \dots \wedge C_m$ where each C_i (called a *clause*) is a disjunction of literals in the set $\{x_1, \bar{x}_1, x_2, \bar{x}_2, \dots, x_n, \bar{x}_n\}$. We want to find an assignment of the Boolean variables that satisfies as many clauses in the formula as possible.

- Write down a randomized algorithm that outputs an assignment with the guarantee that the expected number of clauses satisfied is at least $\frac{m}{2}$.
- Derandomize this algorithm using the method of conditional probabilities discussed in class to get a deterministic algorithm that satisfies at least $\frac{m}{2}$ number of clauses.
- Suppose k is the minimum number of literals in any clause, how will you modify the parameters in part(a) and (b)

2.3 Method of Pessimistic Estimators

We now see a more sophisticated adaptation of the method of conditional expectation. For this we need another randomized algorithm, as it is again going to be algorithm specific⁷.

2.3.1 Congestion Minimization Problem

The problem that we are going to use as the example is called CONGESTION MINIMIZATION problem. We are given a directed graph G with k pairs of vertices $(s_1, t_1), (s_2, t_2), \dots, (s_k, t_k)$ which are sources and destinations respectively. We want to route packets from sources to destination through edge disjoint paths in the graphs so that there is no edge that is used for two paths and hence no change of a congestion. However, even testing whether there are edge-disjoint paths between such pairs of vertices is NP-hard for directed graphs even for $k = 2$. Hence, we have to allow congestion, but ideally we would like to minimise congestion. Formally, a collection of paths P_1, P_2, \dots, P_k (which need not be vertex disjoint) in the above problem is a solution with congestion C , if every edge takes part in at most C paths in the collection. Indeed, we would like to find out a collection of paths with least congestion. The case when $C = 1$ is exactly the edge disjoint path problem which indicates that the optimization problem is hard to solve.

Next step is to look for approximation algorithms. We would like to design algorithm which outputs a set of paths with a guarantess that the congestion is at most αC where C is the optimal congestion possible for the given graph and the $\{(s_i, t_i)\}_{1 \leq i \leq k}$ pairs. Indeed, the problem admits a randomized approximation algorithm with a reasonable approximation ratio.

⁷That is, it can be applied for a class of algorithms for the problem which has the property.

THEOREM 2.3.1. CONGESTION MINIMIZATION admits a randomized algorithm where the solution is guaranteed to be at most $\left(\frac{\log n}{\log \log n}\right) \times \text{OPT}$ where OPT is the optimal congestion possible for the input instance and n is the size of the graph.

We need some details about this randomized algorithm and the analysis since we have to deal with specifics of the analysis in the method itself.

2.3.2 Randomized Approximation Algorithm

The algorithm uses a standard and quite effective technique of designing randomized algorithms which is *linear programming relaxation and randomized rounding*. Firstly the problem can be written as a linear program as follows:

LP Formulation: For any $i \in [k]$, let \mathcal{P}_i denote the set of paths in the graph G which are from vertices s_i to t_i . In the solution, exactly one of these paths needs to be chosen. Let P be the notation for one such path. Let us introduce a Boolean variable x_i^P to indicate whether the path $P \in \mathcal{P}_i$ is chosen for the solution or not. The following k constraints say, that from each \mathcal{P}_i exactly one path must be chosen for the solution :

$$\forall i \in [k], \sum_{P \in \mathcal{P}_i} x_i^P = 1 \quad (2.4)$$

And, we need to minimise the congestion bound C subject to above constraints. Writing down this objective function mathematically:

$$C = \sum_{e \in E} \sum_{\substack{P \in \mathcal{P}_i \\ e \in P}} x_i^P \quad (2.5)$$

Unfortunately, solving such "integer linear programs" is NP-hard. But a technique is that of linear programming relaxation, where we do not insist anymore that the variables must take Boolean values. But instead, we allow them to be real numbers under the constraint - $\forall i, P, 0 \leq x_i^P \leq 1$ which are again linear constraints. There are standard techniques by which this can be solved now in polynomial in the number of variables. However, note that the number of variables in our setting is exponential in n . But we will ignore this fact⁸ for now as the exposition of the technique is easier with the above set up.

Randomized Rounding: The solution to the above linear program gives us values for x_i^P (call them α_i^P between 0 and 1) for $i \in [k]$ and $P \in \mathcal{P}_i$ such that the congestion expression is at most C^* (and is optimal). But this does not point to choice of any path $P \in \mathcal{P}_i$ as they are not Boolean variables. So we need to make that choice (and hence turn the values to Boolean) and this is the place where randomness is used⁹.

⁸The trick to address this issue is to formulate a linear program with the granularity of edges - that is introducing variables at the edge-level than path-level as we have done.

⁹A simple deterministic idea is worth thinking about - namely, for every i , choose the P for which α_i^P is maximum to be the path - this amounts to assigning x_i^P to one for that path P but 0 for the other paths. But this can increase the congestion (constraint ??) without any reasonable bound.

The idea is as follows. For each i , we will choose a path $P \in \mathcal{P}_i$ at random with probability x_i^P . Notice that this is a probability distribution because of constraint 2.4. This can also be seen as, for each i , we will choose one of the x_i^P at random with probability α_i^P and then assign that variable to 1 and make the rest of the variables for that i to be 0. This is equivalent to choosing one path $P_i \in \mathcal{P}_i$ to be in the solution. For example, say we have three paths with values 0.2, 0.7, and 0.1. Get a random number between 0 and 1. If the number is between 0 and 0.2, pick the first path. If the number is between 0.2 and 0.9, pick the second path, and if the number is between 0.9 and 1, pick the third path.

Bounding the Approximation: We need to analyse how badly the objective function (which was evaluating to optimal value C when α_i^P were the assignments for x_i^P s) be affected by this *rounding* step. For every edge $e \in E$, define a random variable Y_i^e as follows:

$$Y_i^e = \begin{cases} 1 & \text{if } e \in P_i \\ 0 & \text{otherwise} \end{cases} \quad \text{and by definition, } \mathbb{E}[Y_i^e] \leq \Pr[e \in P_i] \leq \sum_{\substack{P \in \mathcal{P}_i \\ e \in P}} \alpha_i^P$$

Define $Y_e = \sum_i Y_i^e$. By linearity of expectation, $\mathbb{E}[Y_e] = \sum_i \mathbb{E}[Y_i^e] = C^*$ (because it is exactly the objective function). Now we are in a perfect situation for a tail bound. What is the probability that the random variable (which in this case is a sum of independent random variables).

PROPOSITION 2.3.2 (Chernoff Bound). *Let $X = \sum X_i$ be a random variable expressed as a sum of X_i 's which are independent random variables. Let $\mathbb{E}[X] \leq \mu$. Then for any $\alpha > 1$,*

$$\Pr[X \geq \alpha\mu] \leq e^{-\mu \times [\alpha \ln \alpha - \alpha + 1]}$$

We just need to apply this with a $\alpha = \frac{\log n}{\log \log n}$ is the approximation ratio that we are looking for. Working out the math, this gives $\Pr[Y_e > \alpha C^*] \leq \frac{1}{n^3}$. By union bound, probability that there exists an edge $e \in E$ with $Y_e > \alpha C^*$ is at most $\sum_e \Pr[Y_e > \alpha C^*] \leq \frac{1}{n}$. In other words, with probability at least $1 - \frac{1}{n}$ the above rounding gives a solution (equivalently a Boolean assignment for the variables) which satisfies all the constraints but at the same time does not make the objective function value worse than a factor of α . Hence it is a randomized alpha approximation algorithm.

2.3.3 Pessimistic Estimator

We abstract out a scenario from the above analysis as follows. Let $X = X_1 + X_2 + \dots + X_k$ be a random variable that is expressed as a sum of independent random variables in $[0, 1]$. Chernoff Bound is typically applied when we want to estimate $\Pr[\sum_i X_i > \alpha]$. To see the event as a Boolean value, define the indicator function:

$$I(x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } \sum_i x_i \geq \alpha \\ 0 & \text{otherwise} \end{cases}$$

So one way to descend down the tree to derandomize is to apply the method of conditional expectation on $I(x_1, x_2, \dots, x_n)$. Similar to evaluating $e(r_1, r_2, \dots, r_i)$ in the MAXCUT algorithm, we should

be able to efficiently evaluate:

$$\mathbb{E}[I(x_1, x_2, \dots, x_i, X_{i+1}, \dots, X_n)] = \Pr \left[\sum_{j=1}^k X_j \geq \alpha \mid X_1 = x_1, X_2 = x_2, \dots, X_i = x_i \right]$$

But unfortunately, the RHS expression is not easy to compute since it depends on the distribution of the X_j variables. However, we can still work with the framework, with the idea of replacing it with an upper bound function instead. Since we are bounding a bad event of the sum being large, a function which gives a larger value would be more "pessimistic" and hence the term *pessimistic estimator*.

We can use an exponentiation idea coming from Chernoff Bound proof: For any $t > 0$:

$$\Pr \left[\sum_{j=1}^k X_j \geq \alpha \right] \leq \mathbb{E} \left[\frac{\prod_{j=1}^k e^{tX_j}}{e^{t\alpha}} \right]$$

This upper bound is easy to compute when you have substitutions for some of the X_i s. That is,

$$\Pr \left[\sum_{j=1}^k X_j \geq \alpha \mid X_1 = x_1, X_2 = x_2, \dots, X_i = x_i \right] \leq \left(\frac{\prod_{j=1}^i e^{tx_j}}{e^{t\alpha}} \right) \mathbb{E} \left[\prod_{j=i+1}^k e^{tX_j} \right]$$

Since we know the distribution of each X_j , we can compute the $\mathbb{E} \left[\prod_{j=i+1}^k e^{tX_j} \right]$ in the right hand side as we do in the proof of Chernoff Bound proof itself.

Applying the Method to the Algorithm for Congestion Minimization: Recall the random variables involved. For every edge $e \in E$, we used random variable $Y_i^e = \begin{cases} 1 & \text{if } e \in P_i \\ 0 & \text{otherwise} \end{cases}$ Define $Y_e = \sum_i Y_i^e$. By linearity of expectation, $\mathbb{E}[Y_e] = \sum_i \mathbb{E}[Y_i^e] = C^*$. In the analysis of the algorithm, we proved, by using Chernoff bound that, for $\alpha = \frac{\log n}{\log \log n}$:

$$\Pr \left[\sum_i Y_i^e > \alpha C^* \right] \leq \frac{1}{n}$$

As per our plan, rather than working with the indicator $I(x_1, x_2, \dots, x_i)$ we will work with the pessimistic estimator, for edge $e \in E$: $\left(\frac{\prod_{j=1}^i e^{tr_j}}{e^{t\alpha}} \right) \mathbb{E} \left[\prod_{j=i+1}^k e^{tY_j^e} \right]$. We will also incorporate union bound into our estimate to use the pessimistic estimator for the whole event as :

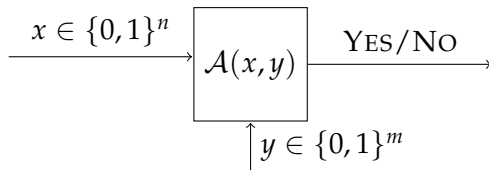
$$I(r_1, r_2, \dots, r_i) \leq f(r_1, r_2, \dots, r_i) = \left(\frac{\prod_{j=1}^i e^{tr_j}}{e^{t\alpha}} \right) \sum_{e \in E} \left(\mathbb{E} \left[\prod_{j=i+1}^k e^{tY_j^e} \right] \right)$$

By notation, we know that $f(r_1, r_2, \dots, r_n) \leq \frac{1}{n}$.

1: Jayalal says: Todo - A few more lines to be completed here about the precise application of the estimators.

In the last week, we saw algorithm specific techniques of derandomization. More precisely, the derandomization uses the peculiarity of the algorithms in the way they use the random bits and the analysis. In this week, we will go back to the abstract set up where we know nothing about the algorithm other than the resource bounds it uses.

We recall some notations first. A randomized algorithm \mathcal{A} on input x runs in time $t(n)$ (where $n = |x|$) and let $y \in \{0,1\}^{m(n)}$ be the concatenation of the unbiased coin toss experiment that the algorithm does during its execution. Notice that $m(n) \leq t(n)$ (we drop the n when it is not required explicitly). If the algorithm runs in polynomial time $t(n) \leq n^c$ for a constant c independent of n .



The guarantee we have is there is an $\epsilon \in (0, \frac{1}{2}]$.

$$\forall x \in \{0,1\}^n, \Pr_{y \in \{0,1\}^r} [A(x, y) \text{ is correct.}] \geq \frac{1}{2} + \epsilon$$

Imagine that we had a success probability of very close to 1. That is, $\epsilon > \frac{1}{2} - \frac{1}{2^m}$. That is, $\forall x \in \{0,1\}^n : \Pr_{y \in \{0,1\}^m} [A(x, y) \text{ is correct.}] > 1 - \frac{1}{2^m}$. Notice that, now the algorithm does not need to use randomness and can fix y to be any string in $\{0,1\}^r$ and run the algorithm \mathcal{A} and the answer is guaranteed to be correct. (This is because, if there exists at least one $y \in \{0,1\}^m$ for which the algorithm errs then the success probability would have been $\leq 1 - \frac{1}{2^m}$.) Thus we would have derandomized the algorithm efficiently.

But how do we achieve such high success probability? Viewing a randomized algorithm as an experiment that we do in physics lab to compute a value, we will repeat the experiment and take the most frequent value in order to reduce the error. However, this also increases the number of random bits which goes against the above plan. Let us formally review this amplification method nevertheless.

3.1 Success Probability Amplification by Repetition

We first write down the algorithm which follows the simple idea of repetition with independent random bits.

Why would this improve the success probability? and if so, how does it depend on k ? The following lemma answers these. Fix the input x . Let \mathcal{E} represent the event that \mathcal{A} accept on the random string y .

Algorithm 3.6 (\mathcal{A}') : input $x \in \{0,1\}^n$

- 1: $count \leftarrow 0$.
 - 2: Choose k independent random strings $y_1, y_2, \dots, y_k \in \{0,1\}^r$. \triangleright Uses $O(kn)$ random bits.
 - 3: **for each** $i \in [k]$ **do**
 - 4: If $\mathcal{A}(x, y_i)$ accepts, if so increment $count$
 - 5: **end for**
 - 6: If $[count > \frac{k}{2}]$ then output YES else output NO.
-

LEMMA 3.1.1. If \mathcal{E} is an event that $\Pr(\mathcal{E}) \geq \frac{1}{2} + \epsilon$, then the probability the \mathcal{E} occurs atleast $\frac{k}{2}$ times on k independent trials is at least $1 - \frac{1}{2}(1 - 4\epsilon^2)^{\frac{k}{2}}$

Proof. Let q denote the probability the \mathcal{E} occurs atleast $\frac{k}{2}$ times on k independent trials. Let $q_i = \Pr(\mathcal{E} \text{ occurs exactly } i \text{ times in } k \text{ trials})$, $0 \leq i \leq k$. Thus, $q = 1 - \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} q_i$. We will analyse the complementary event: $\Pr(\mathcal{E} \text{ occurs atmost } \frac{k}{2} \text{ times}) = \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} q_i$.

We show an upper bound on each q_i and thus show an lower bound on q .

$$\begin{aligned} q_i &= \binom{k}{i} \left(\frac{1}{2} + \epsilon\right)^i \left(\frac{1}{2} - \epsilon\right)^{k-i} \\ &\leq \binom{k}{i} \left(\frac{1}{2} + \epsilon\right)^i \left(\frac{1}{2} - \epsilon\right)^{k-i} \left(\frac{\frac{1}{2} + \epsilon}{\frac{1}{2} - \epsilon}\right)^{\frac{k}{2}-i} \quad (\text{because } \epsilon \leq \frac{1}{2}) \\ &= \binom{k}{i} \left(\frac{1}{2} + \epsilon\right)^{\frac{k}{2}} \left(\frac{1}{2} - \epsilon\right)^{\frac{k}{2}} \\ &= \binom{k}{i} \left(\frac{1}{4} - \epsilon^2\right)^{\frac{k}{2}} \end{aligned}$$

Now we analyse the sum:

$$\begin{aligned} \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} q_i &\leq \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} \binom{k}{i} \left(\frac{1}{4} - \epsilon^2\right)^{\frac{k}{2}} \\ q = 1 - \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} q_i &\geq \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} \binom{k}{i} \left(\frac{1}{4} - \epsilon^2\right)^{\frac{k}{2}} \\ &= 1 - \left(\frac{1}{4} - \epsilon^2\right)^{\frac{k}{2}} 2^{k-1} \\ &= 1 - \frac{1}{2} (1 - 4\epsilon^2)^{\frac{k}{2}} \\ \text{Thus, } q &\geq 1 - \frac{1}{2} (1 - 4\epsilon^2)^{\frac{k}{2}} \end{aligned}$$

□

Thus, if we had an algorithm \mathcal{A} with $\epsilon = \frac{1}{3}$ (that is success probability is at least $\frac{2}{3}$, and we want an algorithm with \mathcal{A}' with $\epsilon = \frac{1}{4}$ (that is, success probability is at least $\frac{3}{4}$). Then, the number of

times the iteration that needs to be done can be back calculated as the k that satisfies:

$$1 - \frac{1}{2} \left(1 - \frac{4}{9}\right)^{\frac{k}{2}} \geq \frac{3}{4}$$

which will be a constant. Quite interestingly, we can do this even when the required error probability is exponentially small. That is, suppose we require the success probability to be $1 - \frac{1}{2^{q(n)}}$ which is quite close to 1. Then the value of k should be :

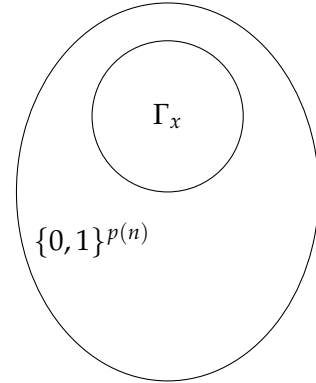
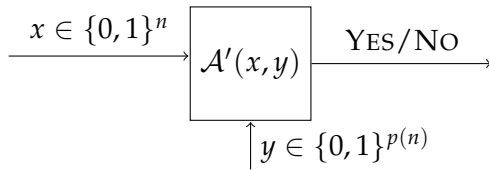
$$1 - \frac{1}{2} \left(1 - \frac{4}{9}\right)^{\frac{k}{2}} \geq 1 - \frac{1}{2^{q(n)}} \implies \left(\frac{9}{5}\right)^{\frac{k}{2}} \leq 2^{q(n)-1}$$

which in turn would imply $k = p(n)$ to be a polynomial value in terms of n . Thus we have the following lemma:

LEMMA 3.1.2 (Amplification Lemma). *Let \mathcal{A} be a randomized algorithm running in time $\text{poly}(n)$ which has $\frac{1}{2} + \epsilon$ as the probability of success. Then for any $q(n)$, we have a randomized algorithm \mathcal{A}' that runs in time $\text{poly}(n)$ with the following success probability. For every input $x \in \{0, 1\}^n$.*

$$\Pr_y [\mathcal{A}'(x, y) \text{ is correct}] \geq 1 - 2^{-q(n)}$$

Define $\Gamma_x = \{y \in \{0, 1\}^{p(n)} \mid \mathcal{A}(x, y) \text{ is correct}\}$.



The guarantee we have is : $\forall x, |\Gamma_x| \geq (1 - 2^{-q(n)})2^{p(n)}$

Exercise 3.1.3 (See Problem Set 1(Problem 5)). A decision problem L is in a class called BPP if there exists a randomized polynomial-time algorithm A such that for every $x \in L$ it holds that $\Pr[A(x, y) = 1] \geq \frac{2}{3}$, and for every $x \notin L$ it holds that $\Pr[A(x) = 0] \geq \frac{2}{3}$. For $\epsilon : \mathbb{N} \rightarrow [0, 1]$, let BPP_ϵ denote the class of decision problems that can be solved in probabilistic polynomial time with error probability upper-bounded by ϵ . Prove the following two claims:

- (a) For every positive polynomial p and $\epsilon(n) = \frac{1}{2} - \frac{1}{p(n)}$, the class BPP_ϵ equals BPP.
- (b) For every positive polynomial p and $\epsilon(n) = 2^{-p(n)}$, the class BPP_ϵ equals BPP.

We already proved something similar in class (See Amplification Lemma). This exercise asks you to prove the same using tail bounds. Given an algorithm A , consider an algorithm A' that on input x invokes A on x for $t(|x|)$ times, and decided based on majority as we did in class. For Part (a) $t(n) = O(p(n)^2)$ and apply Chebyshev's Inequality. For Part 2 set $t(n) = O(p(n))$ and apply the Chernoff Bound.

3.2 Sipser's Argument

From the previous section, we concluded that, for every $x \in \{0,1\}^{p(n)}$, there is a large set Γ_x of random strings which are "good" for x . However, we do not know how to find even an element $y \in \Gamma_x$ in time $\text{poly}(n)$. If we do, then we have a derandomization for the algorithm \mathcal{A}' (which derandomizes \mathcal{A} too).

It is intuitive to think that since a large fraction of $y \in \{0,1\}^{p(n)}$ are good for each $x \in \{0,1\}^n$, could there be a single y which is good for all x ? We show that this is indeed the case. This, in fact, is a simple consequence of the amplification lemma in the previous section.

We first apply Lemma 3.1.2 with $q(n) = 2n$. Thus we have a randomized polynomial time algorithm with error bound 2^{-2n} (i.e. at most 2^{-2n} fraction of random strings are "bad") using $p(n)$ randomness. More precisely, we have that:

$$\text{Answer is YES for input } x \Rightarrow \Pr_y [\mathcal{A}(x, y) \text{ accepts}] \geq 1 - 2^{-2n} \quad (3.6)$$

$$\text{Answer is NO for input } x \Rightarrow \Pr_y [\mathcal{A}(x, y) \text{ accepts}] \geq 2^{-2n} \quad (3.7)$$

That is in such a machine the number of random strings y which lead the machine to output a wrong answer is bounded by 2^{-2n} . Now let us consider a matrix M whose rows are indexed by inputs of length n and columns are indexed by random strings of length $p(n)$, and the (x, y) th entry is 1 if on fixing the random bits to be y the machine M on input x outputs correctly and it is 0 otherwise. That is $M[x, y] = 1$ if and only if $\mathcal{A}'(x, y)$ is correct for the input x . By the amplification we are guaranteed that for a given input x at most 2^{-2n} fraction of the random strings can have $M[x, y] = 0$. Hence the total number of zeros in the A matrix is at most the number of rows times the maximum number of zeros in a row, which is equal to

$$\begin{aligned} \# \text{0's in matrix } M &\leq 2^n \times 2^{-2n} \times 2^{p(n)} \\ &\leq 2^{p(n)-n} \end{aligned}$$

But the total number of zeros, $2^{p(n)-n}$ is strictly less than the number of columns in the matrix M . Hence there must be at least one column with no zeros in it. If a column in the M matrix has no zeros then by the definition of M matrix, the random string w represented by this column when fed as random bits to machine $\mathcal{A}'(x, w)$ would output the correct answer for every $x \in \{0,1\}^n$. Thus we have the following lemma.

LEMMA 3.2.1. *For every n , there is a $y \in \{0,1\}^{p(n)}$ such that $y \in \Gamma_x$ for every $x \in \{0,1\}^n$.*

Thus there is a function $h : \mathbb{N} \rightarrow \Sigma^*$ such that the answer to x is YES if we run the algorithm \mathcal{A} on input x and random string $h(|x|)$ it is guaranteed not to err. This will give a complete derandomization. However, we need the function h to be computable in polynomial time and that is a challenge.

3.3 Amplification with Dependent Trials using Expanders

We study another seemingly unrelated graph theoretic object which are called expanders. They have been extensively studied and used in many areas of science and engineering as the set of sparse graph families which has high connectivity properties. These are particularly desirable in network design such that the network is highly fault tolerant - even the failure of a few links (edges) will not affect the connectivity of the network. The fact that this can be achieved without having too many edges (complete graph is a trivial way to do this) is what makes expander graphs more applicable in this setting.

Expanders have been studied in the theoretical side as well for past few decades and have found numerous applications. Informally, the graph is such that every subset of vertices fetches a large set of vertices in its immediate neighborhood. For a set of vertices in a graph G , define $N(S)$ to be the neighbors of S . We define the graph class formally now.

DEFINITION 3.3.1. (Expander Graphs) A graph $G(V, E)$ is said to be a (α, β) -expander if every $S \subseteq V$ such that $|S| \leq \alpha|V|$, the number of new neighbors $|N(S) \setminus S| \geq \beta|S|$.

A trivial (but unfortunately useless) example of an expander graph is the complete graph. Indeed, in a complete graph on n vertices, if we choose any $S \subseteq V$, such that $|S| \leq \frac{n}{2}$, we have that $N(S) \setminus S = V \setminus S$. Hence, $|N(S) \setminus S| \geq |S|$. This gives that it is a $(\frac{1}{2}, 1)$ -expander as per the above definition. However, it is going to be useless as we will see in our application. In fact, ideally, we would like expander graphs where the degree d and α, β are all constants independent of n .

Thus, we are looking for graphs with constant degree (if possible, even regular). It implies that the graph must be sparse¹⁰ Let us also record a non-example of an expander, which is even a sparse graph. Consider an $n \times n$ complete grid graph, which has n^2 vertices and all possible *grid edges* present. Consider a set S which is of size \sqrt{n} , which has n vertices and hence $|S| \leq \alpha|V|$ for any constant α . However, the number of new neighbors for the set in the grid graph will be $O(\sqrt{n})$ which is not at most $\beta|S|$ for any constant β . Hence such a graph is not an expander for any constant α and β .

To ensure that the definition of expander is practiced enough, we suggest the following exercise which derives a combinatorial consequence of the expansion property.

Exercise 3.3.2 (See Problem Set 1(Problem 6)). Let $G(V, E)$ be an undirected graph n vertices which is $(\frac{1}{2}, 2)$ -expander. Show that the diameter of the graph is at most $O(\log n)$. The diameter of a graph is at most k if and only if between any two vertices in the graph G , there is a path of length at most k in the graph G .

Random Walks on Expanders and Our Application: Although the above definition explains the name expander graphs, the application in our context comes from the fact that if we choose a vertex in the expander at random and then choose the next vertices uniformly at random, the distribution of the ℓ^{th} vertex that you get to (as ℓ increases) is very close to uniform over the entire connected component of the graph where your starting vertex was lying. Intuitively, this is termed as the *rapid mixing* of random walks on expander walks.

¹⁰a graph is said to be sparse if $|E| = O(|V|)$, or else it is called dense.

We will prove the technical details of this idea at a later point in the course where we introduce expander graphs. Informally, to apply the expander graphs in the context of the success probability amplification, we consider the graph on G as a graph on $2^{p(n)}$ vertices indexed by $\{0, 1\}^{p(n)}$. That is, each vertex in the graph can be interpreted as a $y \in \{0, 1\}^{p(n)}$. Based on this, the following algorithm gives a good amplification of success probability.

Algorithm 3.7 (\mathcal{A}') : input $x \in \{0, 1\}^n$

- 1: $count \leftarrow 0$.
 - 2: Let $G(V, E)$ be an expander graph on $2^{p(n)}$ vertices.
 - 3: Choose a vertex $y_1 \in V$ uniformly at random. \triangleright Uses $O(p(n))$ random bits
 - 4: Starting at v perform a random walk for k steps in G . Let $y_1, y_2, \dots, y_k \in \{0, 1\}^r$ be the vertices representing the walk. \triangleright Uses $O(k \log d)$ random bits.
 - 5: **for each** $i \in [k]$ **do**
 - 6: If $\mathcal{A}(x, y_i)$ accepts, if so increment $count$
 - 7: **end for**
 - 8: If $[count > \frac{k}{2}]$ then output YES else output NO.
-

The idea of the above algorithm is based on the rapid mixing of random walks. Without the details, the strings y_1, y_2, \dots, y_k are “almost” as good as randomly chosen y_i ’s since the random walk mixes fast and hence the amplification (although with weaker parameters) can be derived. This will be done in upcoming lectures.

Notice that the number of random bits used by the algorithm is $O(k \log d)$. This explains why the complete graph (whose degree in this case would have been $2^{p(n)} - 1$) would not have been good enough for our purpose.

Part II

Exercise & Problem Sets

Chapter 4

Exercises

4.1 Exercises

Exercise 4.1.1. Let $G \in G(n, p)$. For all $S \subseteq V$, let A_S be the event that S forms an independent set in G . Show that if S and T are two distinct subsets of k vertices then A_S and A_T are independent if and only if $|S \cap T| \leq 1$.

Exercise 4.1.2. Prove that an event \mathcal{E} is independent of a set of events $\{\mathcal{E}_j \mid j \in J\}$ if and only if for all $J_1, J_2 \subseteq J$ such that $J_1 \cap J_2 = \emptyset$

$$\Pr[\mathcal{E} \cap (\cap_{j \in J_1} B_j) \cap (\cap_{j \in J_2} \overline{B_j})] = \Pr(\mathcal{E}) \Pr[(\cap_{j \in J_1} B_j) \cap (\cap_{j \in J_2} \overline{B_j})]$$

4.2 Curiosity Drive

Here we list down all the "out of curious" questions that we discussed (sometimes even not discussed) in the class (and hence in this document).

Chapter 5

Problem Sets

5.1 Problem Set #1

- (1) (See Exercise 2.1.6) A random k -colouring for a graph G is an element of the probability space (Ω, Pr) where Ω is the set of all k -colourings (i.e. partition of V into k sets (V_1, V_2, \dots, V_k) , all this colourings being equally likely (so happening with probability $\frac{1}{k^n}$). For every edge e of G , let A_e be the event that the two endvertices of e receive the same colour. Show that:
 - (a) for any two edges e and f of G , the events A_e and A_f are independent.
 - (b) if e, f and g are three edges of a triangle of G , the events A_e, A_f and A_g are dependent.
- (2) (See Exercise 2.1.7) A graph $G = (V, E)$ is created at random by selecting each edge with probability p . What is the expected number of spanning trees in the randomly sampled graph? (Hint : Use Cayley's Theorem that the number of distinct spanning trees on n vertices is n^{n-2} . Order them, and define an indicator random variable.)
- (3) (See Exercise 2.2.1) In the derandomization of MAXCUT algorithm that we described, we derived an expression for $V(r_1, r_2, \dots, r_i)$ for any i and $r_1, r_2, \dots, r_i \in \{0, 1\}$. We used this to determine, the value of r_{i+1} by computing $V(r_1, r_2, \dots, r_i, 0)$ and $V(r_1, r_2, \dots, r_i, 1)$ and then choosing the value of r_{i+1} to be the one which produces the largest among the two. Prove that the choice of r_{i+1} will be 1 if vertex $i + 1$ has more neighbors in T_i than in S_i and vice versa. Hence, write down the derandomized 0.5-approximation deterministic polynomial time algorithm for MAXCUT as a simple greedy algorithm in terms of the above rule.
- (4) (See Exercise 2.2.2) Let $x_1, x_2, \dots, x_n \in \{0, 1\}$ be Boolean variables and let f be a Boolean formula in CNF form. That is, $f = C_1 \wedge C_2 \wedge \dots \wedge C_m$ where each C_i (called a *clause*) is a disjunction of literals in the set $\{x_1, \bar{x}_1, x_2, \bar{x}_2, \dots, x_n, \bar{x}_n\}$. We want to find an assignment of the Boolean variables that satisfies as many clauses in the formula as possible.
 - (a) Write down a randomized algorithm that outputs an assignment with the guarantee that the expected number of clauses satisfied is at least $\frac{m}{2}$.
 - (b) Derandomize this algorithm using the method of conditional probabilities discussed in class to get a deterministic algorithm that satisfies at least $\frac{m}{2}$ number of clauses.

- (c) Suppose k is the minimum number of literals in any clause, how will you modify the parameters in part(a) and (b)
- (5) (See Exercise 3.1.3) A decision problem L is in a class called BPP if there exists a randomized polynomial-time algorithm A such that for every $x \in L$ it holds that $\Pr[A(x, y) = 1] \geq \frac{2}{3}$, and for every $x \notin L$ it holds that $\Pr[A(x) = 0] \geq \frac{2}{3}$. For $\epsilon : \mathbb{N} \rightarrow [0, 1]$, let BPP_ϵ denote the class of decision problems that can be solved in probabilistic polynomial time with error probability upper-bounded by ϵ . Prove the following two claims:
- (a) For every positive polynomial p and $\epsilon(n) = \frac{1}{2} - \frac{1}{p(n)}$, the class BPP_ϵ equals BPP.
- (b) For every positive polynomial p and $\epsilon(n) = 2^{-p(n)}$, the class BPP_ϵ equals BPP.

We already proved something similar in class (See Amplification Lemma). This exercise asks you to prove the same using tail bounds. Given an algorithm A , consider an algorithm A' that on input x invokes A on x for $t(|x|)$ times, and decided based on majority as we did in class. For Part (a) $t(n) = O(p(n)^2)$ and apply Chebyshev's Inequality. For Part 2 set $t(n) = O(p(n))$ and apply the Chernoff Bound.

- (6) (See Exercise 3.3.2) Let $G(V, E)$ be an undirected graph n vertices which is $(\frac{1}{2}, 2)$ -expander. Show that the diameter of the graph is at most $O(\log n)$. The diameter of a graph is at most k if and only if between any two vertices in the graph G , there is a path of length at most k in the graph G .