# Solving Random Satisfiable 3CNF Formulas in Expected Polynomial Time

MICHAEL KRIVELEVICH [*]        DAN VILENCHIK [†]

October 9, 2005

## Abstract

We present an algorithm for solving 3SAT instances. Several algorithms have been proved to work **whp** (with high probability) for various SAT distributions. However, an algorithm that works **whp** has a drawback. Indeed for typical instances it works well, however for some rare inputs it does not provide a solution at all. Alternatively, one could require that the algorithm always produce a correct answer but perform well on average. *Expected polynomial time* formalizes this notion. We prove that for some natural distribution on 3CNF formulas, called planted 3SAT, our algorithm has *expected polynomial* (in fact, almost linear) running time. The planted 3SAT distribution is the set of satisfiable 3CNF formulas generated in the following manner. First, a truth assignment is picked uniformly at random. Then, each clause satisfied by it is included in the formula with probability $p$. Extending previous work for the planted 3SAT distribution, we present, for the first time for a satisfiable SAT distribution, an *expected polynomial time* algorithm. Namely, it solves all 3SAT instances, and over the planted distribution (with $p = d/n^2$, $d > 0$ a sufficiently large constant) it runs in expected polynomial time. Our results extend to $k$-SAT for any constant $k$.

# 1 Introduction and Results

## 1.1 The Planted 3SAT Distribution

A 3CNF formula over the variables $x_1, x_2, ..., x_n$ is a conjunction of clauses $C_1, C_2, ..., C_m$ where each clause is a disjunction of three literals. Each literal is either a variable or its negation. A 3CNF is satisfiable if there is a boolean assignment to the variables s.t. every clause contains at least one literal which evaluates to true. 3SAT is the language of all satisfiable 3CNF formulas.

Although 2SAT is known to be in P, 3SAT is one of the most famous NP-complete problems [Coo71]. Understanding the complexity of the satisfiability problem on various natural instances is a fundamental problem in computer science. Thus, while the problem is believed unsolvable in

---

polynomial time in the worst case, many heuristics have been proposed and some of them seem to perform quite well "on average".

Algorithmic theory of random structures has been the focus of extensive research in recent years (see [FM97] for a comprehensive survey). As part of this trend, uniformly random 3CNFs (generated by selecting at random $m = m(n)$ clauses over the variables $\{x_1, ..., x_n\}$) have also caught the eyes of many researchers. Random 3SAT is known to have a sharp satisfiability threshold in the clause-to-variable ratio [Fri99]. Namely, a random 3CNF with clause-to-variable ratio below the threshold is satisfiable **whp** (with high probability) and one with ratio above the threshold is unsatisfiable **whp**. This threshold is not known exactly (and not even known to be independent of $n$). The threshold is known to be at least 3.42 [KKL02] and at most 4.5 [KKS$^+$01]. Experimental results predict the higher end of the interval [CA96]. In this work we concentrate on formulas with constant (above the threshold) clause-to-variables ratio. At such ratios, most formulas are unsatisfiable, and analyzing the distribution of random satisfiable instances seems hard. Thus, we focus on the *planted* 3SAT distribution, denoted throughout by $\mathbb{P}_{p,n}$.

A planted instance in $\mathbb{P}_{p,n}$ is generated by first picking at random a truth assignment $\varphi$ to the $n$ variables. Next, each clause satisfied by $\varphi$ is included with probability $p$. In this work we consider $p \geq d/n^2$, where $d$ is a sufficiently large constant. The distribution is highly interesting in its own right. It has been studied empirically in [CI95], and is the analog of the planted clique, planted bisection, and planted coloring distributions studied e.g. in [AK97], [AKS98], [BS95], [FK00]. The planted 3SAT distribution is studied in [KP92], where an algorithm based on greedy variable assignment rule is proved to work **whp** for dense instances ($p = \Theta(1/n)$). The authors conjecture that some modification of their algorithm will work for sparser instances. [Fla03], drawing on techniques from [AK97] (spectral and combinatorial techniques), proves this conjecture.

## 1.2 Expected Polynomial Time

An algorithm that works **whp** has, nevertheless, a drawback. Indeed for typical instances it works well, however for some rare inputs it does not provide a solution at all. Alternatively, one could require that the algorithm always produce a correct answer but perform well on average. Formally, an algorithm $\mathcal{A}$ with running time $t_{\mathcal{A}}(I)$ on an input instance $I$, has *expected polynomial running time* over a distribution $\mathcal{D}$ on the inputs, if $\sum_I t_{\mathcal{A}}(I) \cdot Pr_{\mathcal{D}}[I]$ is polynomial. Note that an expected polynomial time algorithm also works **whp**, the converse however isn't necessarily true. In many cases, in order for an algorithm to run in expected polynomial time, a distinction between typical instances (which are easy, and require polynomial time) and atypical ones (which are rare, but may require superpolynomial time) is due. Moreover, loosely speaking, the amount of work put in should be inversely proportional to how atypical the instance is. To perform this dichotomy, a deeper understanding of the underlying probability space is required, as well as a more careful analysis of properties of the random instance. Therefore, by designing algorithms that work in expected polynomial time, one might expect to obtain more robust and efficient algorithms. Several papers discussed algorithms with expected polynomial time, we selectively mention some: [BV04], [BW85], [Böt05], [CO04], [CO03], [COGLS04], [COT03], [DF89], [FSVM93], [Kri02], [KV02], [PS92], [Wil84].

### 1.3 Related Work

The coloring algorithm in [AK97] works **whp** for sparse planted $k$-colorable graphs, however it is not an expected polynomial time algorithm. An expected polynomial time algorithm for coloring planted $k$-colorable graphs is given in [CO04] for relatively dense graphs (average degree $np = d \cdot k \cdot \log n$). [Böt05] improves upon this result, giving an expected polynomial time algorithm for sparse $k$-colorable graphs (average degree $np = d \cdot k^2$, which is constant when $k$ is). Analogously for the planted 3SAT setting, the algorithm presented in [Fla03] is not an expected polynomial time algorithm. Although expected polynomial time is not discussed in [KP92], slightly changing the algorithm makes it such. However, [KP92] works for dense formulas, namely when $p = d/n$.

### 1.4 Our Results

Combining ideas from [Böt05] and [Fla03], we devise an algorithm for solving (i.e. finding a satisfying assignment) satisfiable 3CNF instances. We prove that our algorithm has expected polynomial running time over $\mathbb{P}_{d/n^2,n}$. As far as the authors know, this is the first work to address the issue of expected polynomial time algorithms for satisfiable SAT distributions. Formally we prove,

**Theorem 1.** *There exists an algorithm SAT that solves 3SAT, and runs in expected polynomial time over the planted 3SAT distribution with $p \geq d/n^2$, $d$ a sufficiently large constant.*

The theorem extends to any constant $k$, while in this case $d$ depends on $k$.

The remainder of the paper is structured as follows. In Section 2, the motivation behind the algorithm, and the algorithm itself are given. In Section 3, a few properties of $\mathbb{P}_{d/n^2,n}$ are analyzed. These properties come useful when estimating the expected running time of the algorithm in Section 4. Finally, in Section 5, we discuss some possibly interesting issues for further research.

## 2 The Algorithm

Before presenting the algorithm, we offer some intuition. The following discussion assumes the input is sampled according to $\mathbb{P}_{d/n^2,n}$.

### 2.1 Motivation

The algorithm looks for the planted assignment (though it may come across another satisfying assignment). Thus, we say that a variable is wrongly assigned by some truth assignment $\tau$, if $\tau(x)$ differs from the planted assignment of $x$.

Suppose that the planted assignment sets $x$ to be true. In every chosen clause that contains $x$, $x$ appears positively (namely, as the literal $x$) with probability 4/7, and negatively (namely, as $\bar{x}$) with probability 3/7. Therefore, by comparing for every variable the number of positive and negative appearances, and assigning it according to the majority, one expects to already have a good approximation of the planted assignment. Such a procedure is often called a *majority vote* (line 1 in the algorithm SAT). In fact, when $p \geq d \log n/n^2$, for a sufficiently large constant $d$, the

majority vote suffices to reconstruct the planted assignment **whp**, see [BSBG02] and [Gen98] for a complete discussion and proof. However, when $p = d/n^2$, the majority vote is wrong for some small fraction of the variables ($\exp\{-\Theta(n^2 p)\}$, which is constant in our case).

The next step is to distinguish between correctly and wrongly assigned variables. This is achieved by a careful unassignment procedure, in which the assignment of variables suspected to be wrong is peeled off, leaving us with a partial assignment (lines 3-5). Now we are in a better position as this partial assignment coincides **whp** with the planted one.

It remains to complete the assignment of the unassigned variables. Consider the formula induced by the unassigned variables. One can prove that this formula breaks down **whp** to mutually disjoint subformulas, each containing at most $O(\log n)$ variables. Therefore, using exhaustive search, separately in every subformula, the partial assignment is completed, if possible, to a satisfying one (line 7,8 and 13), while spending polynomial time.

This 3-step heuristic is very similar to the one described and analyzed in [Vil04] (based on [Fla03]). The heuristic finds **whp** a satisfying assignment for $\mathbb{P}_{d/n^2,n}$ but does not necessarily work for all instances. It might be the case, although rare, that the unassignment step didn't filter out the wrongly assigned variables. Therefore the exhaustive search might fail (as the formula induced by the unassigned variables may no longer be satisfiable, or the partial assignment may induce a $(F \vee F \vee F)$ clause). Since the algorithm is required to produce an answer for all instances, one needs to take care of this rare event as well. Suppose indeed that the exhaustive search fails. Our goal is to identify the (minimal) set of variables that survived the unassignment s.t. by flipping its current assignment (given by the majority vote), the partial assignment can be completed to a satisfying one. To this end, a controlled recovery step is employed (lines 9-11). Iteratively, one goes over all subsets $Y$ of assigned variables (in ascending order of size). For each subset one tries all possible assignments $\pi(Y)$, looking for the aforementioned set and flip its assignment.

## 2.2  Notation

Let $V = \{x_1, x_2, .., x_n\}$ be the set of all variables. For a set of variables $U \subseteq V$, let $G(U, E)$, the *residual graph*, be the graph defined as follows. $U$ is the set of vertices. For every $x_i, x_j \in U$, $(x_i, x_j) \in E$ if there exists some clause in $I$ containing both $x_i, x_j$ (regardless of their polarity). A *partial truth assignment* is a vector of length $n$ over $\{0, 1, *\}$, where $*$ stands for a variable which is not assigned. A variable $x$ *supports* a clause $C$ with respect to a partial assignment $\psi$, if it is the only variable to satisfy $C$ under $\psi$, and the other two variables are assigned by $\psi$. E.g., if $C = (x \vee y \vee z)$ and $\psi(C) = (T \vee F \vee F)$, then $x$ supports $C$ w.r.t. $\psi$. By *simplify $I$* according to $\psi$, when $\psi$ is a partial assignment, we mean: in every clause substitute every assigned variable with the value given to it by $\psi$. Next, if a clause contains a literal which evaluates to true, remove the clause. Otherwise, remove all literals which evaluate to false (the resulting instance is not necessarily in 3CNF form). Denote by $I|_\psi$ the 3CNF $I$ simplified according to $\psi$. We say that $I|_\psi$ *induces an empty clause* if there exists a clause in which all three literals evaluate to false under $\psi$. For a set of variables $A \subseteq V$, denote by $I[A]$ the set of clauses in which all variables belong to $A$.

The algorithm as presented can be optimized using known procedures such as unit-clause-propagation or the pure-literal rule (which can be performed on $I|_{\psi \circ \pi(Y)}$ before the exhaustive search starts). These optimizations were not included in the description of the algorithm for simplicity of presentation and analysis.

Throughout the paper, $\epsilon$ is used to denote some small constant independent of $d$, say $\epsilon = 10^{-5}$. $I$ denotes a random instance from $\mathbb{P}_{d/n^2,n}$, and $\varphi$ its planted assignment. MAJ is the majority vote assignment of $I$.

---

**Algorithm 1** : SAT(I)

---

1: MAJ $\leftarrow$ the majority vote assignment.
2: Set $i = 0$ and $\psi_0 \leftarrow$ MAJ.
3: **while** there exists a variable $x$ whose support is less than $(1 - \epsilon)d/2$ w.r.t. $\psi_i$ **do**
4:     $\psi_{i+1} \leftarrow \psi_i$ with the variable $x$ unassigned.
5:     $i \leftarrow i + 1$.
6: **end while**
7: Let $\psi$ be the final partial assignment, and $U$ be the set of unassigned variables in $\psi$.
8: Construct the residual graph $G(U, E)$ and find the set of connected components $\{\Gamma_i\}$.
9: **for** $y = 0$ to $|V \setminus U|$ **do**
10:     **for all** $Y \subseteq V \setminus U$, $|Y| = y$ and every assignment $\pi(Y)$ to the variables of $Y$ **do**
11:         Let $\psi \circ \pi(Y)$ be $\pi(x)$ if $x \in Y$, and $\psi(x)$ otherwise.
12:         **if** $I|_{\psi\circ\pi(Y)}$ doesn't induce an empty clause **then**
13:             Independently in every $\Gamma_i$, using exhaustive search try to satisfy $I|_{\psi\circ\pi(Y)}[\Gamma_i]$.
14:             **if** the search succeeded **return** the satisfying assignment **end if**
15:         **end if**
16:     **end for**
17: **end for**
18: **return** $I$ is not satisfiable.

---

The correctness of SAT is evident. If the algorithm returns an assignment (line 14), it satisfies $I$ (by the correctness of the exhaustive search and the condition in line 12). If $I$ is satisfiable, in the worst case, $y = |V \setminus U|$ is reached, all possible assignments to the assigned variables are examined, and at least for the one coinciding with some satisfying assignment, the exhaustive search succeeds. If $I$ is not satisfiable, then at least one of the conditions in lines 12, 14 fails every iteration, line 18 is eventually reached and the algorithm answers "not satisfiable". It thus remains to analyze the expected running time over $\mathbb{P}_{d/n^2,n}$.

## 3  Properties of a Random Instance

In this section, a few properties of $\mathbb{P}_{d/n^2,n}$ are analyzed. These properties come useful when estimating the expected running time of the algorithm. Recall the definition of $\mathbb{P}_{d/n^2,n}$. First a truth assignment $\varphi$ to the variables $V = \{x_1, x_2, ..., x_n\}$ is picked uniformly at random. Next, every clause satisfied by $\varphi$ is included in the formula with probability $p = d/n^2$. Clearly $I$ is satisfiable, since at least $\varphi$ satisfies it. There are $(2^3 - 1) \cdot \binom{n}{3}$ clauses satisfied by $\varphi$. The expected size of $I$ is then $p \cdot 7 \cdot \binom{n}{3} = 7dn/6 + o(n)$. As for the support, every variable $x$ supports $\binom{n-1}{2}$ clauses w.r.t. $\varphi$. The expected support of a variable $x$ in $I$ w.r.t. $\varphi$ is then $p \cdot \binom{n-1}{2} = d/2 + o(1)$. This explains the choice in line 3 of the algorithm.

The following inequality is used to bound binomial coefficients: for $k \in [1, n]$, $\binom{n}{k} \leq (ne/k)^k$. Unless stated otherwise, the basis of the logarithm is the natural basis. In the analysis, both $d$ and $n$ are assumed to be sufficiently large. No attempt was made to optimize the constants in the

analysis. Throughout, we let $\alpha_0 = \exp\{-d/5000\}$, $\alpha_1 = \exp\{-\epsilon^2 d/40\}$.

## 3.1 Majority Vote and Support

**Proposition 2.** *Let $\alpha \in [\alpha_0, \epsilon]$. Let $F_{MAJ}$ be a random variable counting the number of variables in $I$ on which the majority vote disagrees with the planted assignment. Then, $Pr[F_{MAJ} \geq \alpha n] \leq \exp\{-d/4000 \cdot \alpha n\}$.*

There are some delicate technicalities in the proof of this proposition due to dependency issues. The full proof is given in the appendix, here only some intuition is offered. For a variable $x$, the number of clauses it appears in positively, and the ones it appears in negatively are two independent binomially distributed **rv** (random variables) with a gap of about $d/2$ in their expectations. Let $F_x$ be the event $\{MAJ(x) \neq \varphi(x)\}$. Applying standard bounds on the tail of the binomial distribution (e.g. Chernoff bounds), $Pr[F_x] \leq \exp\{-\Theta(d)\}$. Using the union bound,

$$Pr[F_{MAJ} \geq \alpha n] \leq \binom{n}{\alpha n} Pr[F_{x_1} \wedge F_{x_2} \wedge ... \wedge F_{x_{\alpha n}}].$$

The difficulty lies in estimating the latter, as the $F_{x_i}$'s are dependent.

**Proposition 3.** *Let $\alpha \in [\alpha_1, 1]$. Let $X_\epsilon$ be a **rv** counting the number of variables whose support w.r.t. $\varphi$ is less than $(1 - \epsilon)d/2$. Then, $Pr[X_\epsilon \geq \alpha n] \leq \exp\{-\epsilon^2 d/20 \cdot \alpha n\}$*

**Proof.** The support of a variable $x$ is a binomially distributed **rv** with expectation $d/2 + o(1)$. Let $F_x$ be the event that $x$ supports less than $(1 - \epsilon)d/2$ clauses with respect to $\varphi$. Applying the Chernoff bound gives $Pr[F_x] \leq \exp\{-\epsilon^2 d/16\}$. Further, the clauses that $x$ supports are all different than the ones $y \neq x$ supports. Hence, the events $F_x, F_y$ are independent. Therefore, in the case of the support, one can easily perform the calculations sketched in Proposition 2. In particular,

$$Pr[F_{x_1} \wedge F_{x_2} \wedge ... \wedge F_{x_{\alpha n}}] = (Pr[F_x])^{\alpha n}.$$

Using standard calculations, the proposition then follows. ∎

## 3.2 Dense Subformulas

The next property to prove is analogous to a property proved in [AK97] for random graphs. Loosely speaking, [AK97] prove that **whp** a random graph doesn't contain a small induced subgraph with a large average degree. Formally, for 3SAT instances,

**Proposition 4.** *For a set of variables $Y \subseteq V$, denote by $F(Y) \subseteq I$ the set of clauses in $I$ containing at least two variables from $Y$. Let $\rho \in (0, 1]$, and $\alpha \in (0, (\rho/16)^4]$. Then, $Pr[\exists Y \subseteq V, |Y| = \alpha n, \exists F(Y) \subseteq I, |F(Y)| \geq \rho d \cdot \alpha n] \leq \exp\{-(\rho d/4 \cdot \log \frac{1}{\alpha}) \cdot \alpha n\}$*

**Proof.** For fixed $Y$, there are $\binom{\alpha n}{2}(n-2)2^3 \leq 4\alpha^2 n^3$ clauses containing at least two variables from $Y$. Each is included in the formula with probability $d/n^2$. Therefore,

$$Pr[\exists Y \subseteq V, |Y| = \alpha n, \exists F(Y) \subseteq I, |F(Y)| \geq \rho d \cdot \alpha n]$$

$$\leq \binom{n}{\alpha n}\binom{4\alpha^2 n^3}{\rho d \cdot \alpha n}(d/n^2)^{\rho d \cdot \alpha n} \leq \left(\frac{en}{\alpha n}\right)^{\alpha n}\left(\frac{e \cdot 4\alpha^2 n^3}{\rho d \cdot \alpha n}\right)^{\rho d \cdot \alpha n} \cdot \left(\frac{d}{n^2}\right)^{\rho d \cdot \alpha n} \leq \left(\frac{e}{\alpha}\right)^{\alpha n}\left(\frac{12\alpha}{\rho}\right)^{\rho d \cdot \alpha n}$$

$$\leq \exp\{\rho d \cdot \alpha n \cdot (4 - \log\frac{1}{\alpha}/2 + \log\frac{1}{\rho})\} \underbrace{\leq}_{\alpha \leq (\rho/16)^4} \exp\{-(\rho d/4 \cdot \log\frac{1}{\alpha}) \cdot \alpha n\}.$$

■

## 3.3 Expanding Set of Variables

We describe a subset of the variables, denoted throughout by $H$, which plays a crucial role in the analysis. The set $H$ has a good expansion property in the following sense: every variable in $H$ supports at least $(1 - \epsilon)d/2$ clauses in $I[H]$ w.r.t. both the planted assignment and the majority vote. $H$ is defined using the following procedure:

*Let $W$ be the set of variables whose majority vote disagrees with $\varphi$.*
*Let $B$ be the set of variables whose support w.r.t. $\varphi$ is less than $(1 - \epsilon/2)d/2$.*
*Let $H \subset V$ be constructed as follows:*

1. *Let $H_0 = V \setminus (B \cup W)$.*

2. *While there exists a variable $a_i \in H_i$ which supports less than $(1 - \epsilon)d/2$ clauses in $I[H_i]$, define $H_{i+1} = H_i \setminus \{a_i\}$.*

3. *Let $a_m$ be the last variable removed at step 2. Define $H = H_{m+1}$.*

**Proposition 5.** *Let $\bar{H} = V \setminus H$, let $\alpha \in [4\alpha_1, (\epsilon/128)^4]$. Then, $Pr[|\bar{H}| \geq \alpha n] \leq \exp\{-\epsilon^2 d/100 \cdot \alpha n\}$.*

**Proof.** Partition the variables in $\bar{H}$ into variables that belong to $B \cup W$, and variables that were removed in the iterative step, $\bar{H}^{it} = H_0 \setminus H$. If $|\bar{H}| \geq \alpha n$, then at least one of $B \cup W$, $\bar{H}^{it}$ has cardinality at least $\alpha n/2$. Consequently,

$$Pr[|\bar{H}| \geq \alpha n] \leq \underbrace{Pr[|B \cup W| \geq \alpha n/2]}_{(a)} + \underbrace{Pr[|\bar{H}^{it}| \geq \alpha n/2 \mid |B \cup W| \leq \alpha n/2]}_{(b)}.$$

Propositions 2, 3 are used to bound $(a)$,

$$Pr[|B \cup W| \geq \alpha n/2] \leq Pr[|W| \geq \alpha n/4] + Pr[|B| \geq \alpha n/4].$$

Using Proposition 2,
$$Pr[|W| \geq \alpha n/4] \leq \exp\{-d/16,000 \cdot \alpha n\}.$$

Proposition 3 yields
$$Pr[|B| \geq \alpha n/4] \leq \exp\{-\epsilon^2 d/80 \cdot \alpha n\}.$$

7

Combining both,

$$Pr[|B \cup W| \geq \alpha n/2] \leq 2\exp\{-\epsilon^2 d/80 \cdot \alpha n\} \leq \exp\{-\epsilon^2 d/90 \cdot \alpha n\}.$$

To bound $(b)$, observe that every variable that is removed in iteration $i$ of the iterative step (step 2), supports at least $\epsilon/2 \cdot d/2$ clauses in which at least another variable belongs to $\{a_1, a_2, ..., a_{i-1}\} \cup B \cup W$. Consider iteration $\alpha n/2$. Assuming $|B \cup W| \leq \alpha n/2$, by the end of this iteration there exists a set containing at most $\alpha n$ variables, and there are at least $(\epsilon/8) \cdot d \cdot \alpha n$ clauses containing at least two variables from it (note that no clause was counted twice as the supporter of a clause is unique). Plugging $\epsilon/8$ in Proposition 4,

$$Pr[|H^{it}| \geq \alpha n/2 \big| |B \cup W| \leq \alpha n/2] \leq \exp\{-(\epsilon d/32 \cdot \log \frac{1}{\alpha}) \cdot \alpha n\} \leq \exp\{-\epsilon d \cdot \alpha n\}.$$

The last inequality is by the upper bound on $\alpha$. The proposition then follows from combining these results. ■

Bounding the probability of $\bar{H}$ being too large, plays a crucial role in the proof of the following property.

**Proposition 6.** *Let* $\beta \in [\Theta(\log n/(n\epsilon^6 d)), 1]$. *Let* $G(\bar{H}, E)$ *be the residual graph induced by* $\bar{H}$. *The probability there exists a connected component of size* $\beta n$ *in* $G(\bar{H}, E)$ *is at most* $\exp\{-\Omega(\epsilon^6 d) \cdot \beta n\}$.

The proof of this proposition is very similar to the one given in [Fla03], [Vil04]. Here only some intuition is offered. A detailed proof is given in the appendix. Recall the well known $G_{n,p}$ distribution for random graphs, where each edge is included independently of the others with probability $p$. If $p < \frac{1}{n}$ then **whp** the largest connected component in $G$ is of size at most $O(\log n)$ (see e.g. [AS00] for a complete discussion). Now consider a random subgraph of size $\alpha n$ of $G_{n,d/n}$, this subgraph is exactly $G_{\hat{n},\alpha d/\hat{n}}$, where $\hat{n} = \alpha n$. If $\alpha d < 1$, **whp** the largest connected component in $G_{\hat{n},\alpha d/\hat{n}}$ is of size $O(\log \hat{n})$. If on the other hand $\alpha d > 1$, in our setting that would imply that $\bar{H}$ is a large set. This however, happens with small probability (Proposition 5). In our case, the residual graph $G(\bar{H}, E)$ is not a truly random graph, and $\bar{H}$ is not a truly random subset of variables. Hence the proof is more delicate.

## 4  Analysis of the Expected Running Time

To begin with, some intuition is given. As mentioned above, the set $H$ plays a crucial role in the analysis. The majority vote sets $H$ correctly (by the construction of $H$). Also, $H$ survives the unassignment step (due to its expansion properties w.r.t. MAJ). Therefore, by the end of the unassignment step, $I[H]$ is surely satisfied by $\psi$ (the partial assignment defined in line 7 of SAT). It remains to satisfy, if not yet satisfied by $\psi$, $I \setminus I[H]$. Normally, $I[H]$ contains almost all clauses (Proposition 5). Therefore **whp**, $I \setminus I[H]$ breaks down to small-size mutually-disjoint subformulas (Proposition 6), allowing the exhaustive search to run efficiently (more precisely, in expected polynomial time). Unfortunately, the unassignment step does not necessarily screen out exactly the set $H$. There may be additional variables, outside of $H$, that also survived the unassignment step. Normally, their assignment coincides with $\varphi$. Although rare, this may not be the case, possibly preventing the exhaustive search from ending up successfully. Let $Y_0$ be the set of variables which prevent the exhaustive search from ending up successfully. If $Y_0$ is non-empty

(which normally doesn't happen), we shall prove that both the majority vote is wrong for $Y_0$ and $Y_0$ is a dense set in the sense of Proposition 4. Both however, happen with small probability. Carrying out all the calculations, one proves that in expectation, the recovery procedure (lines 9-11, responsible for fixing $Y_0$) consumes constant time, allowing the algorithm to run in expected polynomial time.

Formalizing the above, let $Y_0$ be the set of variables used in the iteration when the algorithm finally finds a satisfying assignment for $I$, and $\pi(Y_0)$ be the assignment given to $Y_0$. We introduce some **rv** that help clarify the description and analysis of the running time: $R_{EXH}$ measures the maximal (over all sets $Y \subseteq V$) running time of the exhaustive search (line 13); $R_{SIMPLIFY}$ measures the time it takes to calculate $I|_{\psi \circ \pi(Y)}$ and $\{I|_{\psi \circ \pi(Y)}[\Gamma_i]\}$ (lines 12-13); $R_{MAIN}$ measures the running time of the majority vote, the unassignment step and the decomposition of $G(U, E)$ to connected components (lines 1-8); $R_{REC}$ (REC for recovery) is the number of times the iteration in lines 9-10 is performed. Finally, let $R_{\mathcal{A}}$ be a **rv** measuring the running time of the algorithm. The probability of all **rv** is taken over the coin flips in the choice of $I$. Using the above notations, $R_{\mathcal{A}} \leq R_{MAIN} + R_{REC} \cdot (R_{EXH} + R_{SIMPLIFY})$. Our goal is to prove that $E[R_{\mathcal{A}}] \leq n^c$ for some constant $c \geq 1$. In particular we shall prove $c = 1 + \Theta(1/d)$.

**Proposition 7.** $E[R_{\mathcal{A}}] \leq E[R_{MAIN}] + E[R_{REC} \cdot R_{EXH}] + E[R_{REC}] \cdot E[R_{SIMPLIFY}]$.

**Proof.** The proposition follows from linearity of expectation and the fact that $R_{SIMPLIFY}$ and $R_{REC}$ are independent **rv** (the simplification's result does depend on the set $Y$, but the time it takes to simplify $I$ can be bounded using a function of $|I|$, which is independent of $R_{REC}$). Unfortunately, $R_{EXH}$ and $R_{REC}$ are dependent. ∎

**Proposition 8.** Both $E[R_{MAIN}] = O(dn)$ and $E[R_{SIMPLIFY}] = O(dn)$.

Let $m$ be the number of clauses in $I$. It is easily seen that $E[R_{SIMPLIFY}], E[R_{MAIN}] \leq O(m \cdot n)$, since the most naive way of performing lines 1-8, and calculating $I|_{\psi \circ \pi(Y)}$, $\{I|_{\psi \circ \pi(Y)}[\Gamma_i]\}$ takes at most $O(m \cdot n)$. To prove $E[R_{SIMPLIFY}], E[R_{MAIN}] = O(dn)$ it suffices to show how to perform these procedures in time $O(m + n)$. Since $E[O(m + n)] = O(dn)$ the proposition follows. This requires some extra, however very standard, work. Full details in the appendix.

**Proposition 9.** $E[R_{REC}] \leq O(1)$.

Before proving this proposition we make the following key observation.

**Proposition 10.** Let $Y_0, \pi(Y_0)$ be as above. If $|Y_0| = y_0 > 0$ then (a) the majority vote is wrong for at least $y_0$ variables (b) there are at least $d \cdot y_0/3$ clauses containing at least two variables from $Y_0$.

**Proof.** To prove (a), suppose by contradiction that MAJ is wrong for less than $y_0$ variables. Let $Y'$ be the set of assigned variables on which MAJ is wrong. Consider the iteration in which $Y = Y'$ and $\pi(Y') = \overline{MAJ}$ (line 10). In this iteration, $\psi \circ \pi(Y')$ agrees with $\varphi$ on all assigned variables. Therefore, the exhaustive search (line 13) succeeds (since $\psi \circ \pi(Y')$ can be at least completed to $\varphi$) and the algorithm stops, contradicting the choice of $Y_0$ ($|Y'| < y_0$, therefore $Y' \neq Y_0$).

To prove (b), observe that by the algorithm $Y_0 \subseteq V \setminus U$. Further, $\pi(Y_0)$ must disagree with MAJ on $Y_0$ (otherwise define $Y'$ to be $Y_0$ minus the variables on which $\pi(Y_0)$ and MAJ agree. The algorithm

9

performs with $Y'$ exactly as it would with $Y_0$, but $|Y'| < y_0$).

Now, if $x \in Y_0$, since $x \in V \setminus U$, by the algorithm, $x$ must support at least $(1 - \epsilon)d/2 > d/3$ clauses in which all variables are assigned. $\pi(Y_0)$ flips the assignment of $x$ (w.r.t. MAJ), causing all the clauses it supports to become not satisfied. Therefore, in every such clause there must be at least another variable $z \in Y_0$ (or else, $I|_{\psi \circ \pi(Y_0)}$ induces an empty clause, and the algorithm could not have ended successfully with $Y = Y_0$). Since every clause is counted once (unique support) property $(b)$ follows. ∎

**Proof.** (Proposition 9) For every $y$ (line 9), the number of iterations performed (line 10) is $\binom{n}{y}2^y$. Hence,

$$E[R_{REC}] \leq \sum_{y=0}^{n} Pr[|Y_0| = y] \cdot \sum_{y'=0}^{y} \binom{n}{y'}2^{y'}.$$

To bound this expression, we decompose it according to the different values of $y = \alpha n$. If $|Y_0| = \alpha n \geq \alpha_0 n$, Propositions 2 and 10.a are employed to bound

$$Pr[|Y_0| = \alpha n] \leq \exp\{-d/4000 \cdot \alpha n\}.$$

Otherwise, Propositions 4 with $\rho = 1/3$ and 10.b give,

$$Pr[|Y_0| = \alpha n] \leq \exp\{-(d/12 \cdot \log \frac{1}{\alpha}) \cdot \alpha n\}.$$

For $y \leq n/2$,

$$\sum_{y'=1}^{y} \binom{n}{y'}2^{y'} \leq y\binom{n}{y}2^y \leq \exp\{y \cdot (\log(n/y) + 3)\}.$$

In any case,

$$\sum_{y'=1}^{y} \binom{n}{y'}2^{y'} \leq 3^n.$$

Finally, we are ready to bound $E[R_{REC}]$:

$$E[R_{REC}] = \sum_{y=0}^{n} Pr[|Y_0| = y] \sum_{y'=0}^{y} \binom{n}{y'}2^{y'} \leq$$

$$\underbrace{O(1)}_{\alpha=0} + \sum_{\alpha=1/n}^{\alpha_0} \exp\{(\log \frac{1}{\alpha} + 3) \cdot \alpha n\} \cdot \underbrace{\exp\{-(d/12 \cdot \log \frac{1}{\alpha}) \cdot \alpha n\}}_{\text{Prop. 4}}$$

$$+ \sum_{\alpha_0}^{\epsilon} \exp\{(\log \frac{1}{\alpha} + 3) \cdot \alpha n\} \cdot \underbrace{\exp\{-d/4000 \cdot \alpha n\}}_{\text{Prop. 2}} + \sum_{\alpha=\epsilon}^{1} 3^n \cdot \underbrace{\exp\{-\epsilon d/4000\}}_{\text{Prop. 2}} \leq$$

$$O(1) + \sum_{\alpha=1/n}^{\alpha_0} \exp\{-(d/20 \cdot \log \frac{1}{\alpha}) \cdot \alpha n\} + \sum_{\alpha=\alpha_0}^{\epsilon} \exp\{(-d/4000 + 3 + d/5000) \cdot \alpha n\}$$

$$+ \sum_{\alpha=\epsilon}^{1} \exp\{-\Omega(d) \cdot n\} \leq O(1) + \sum_{\alpha=1/n}^{\alpha_0} \exp\{-\Omega(d^2) \cdot \alpha n\} + \sum_{\alpha=\alpha_0}^{1} \exp\{-\Omega(d) \cdot \alpha n\}$$

$$\leq O(1) + \sum_{y=1}^{n} 2^{-y} = O(1).$$

■

**Proposition 11.** $E[R_{REC} \cdot R_{EXH}] = n^{1+\Theta(1/d)}$.

Before proving this proposition we make the following observation.

**Lemma 12.** *The set $H$ defined in subsection 3.3 is assigned when lines 9-17 are executed.*

**Proof.** $\varphi$ and MAJ coincide on $H$ (by the construction of $H$). Thus, every variable $x \in H$ supports at least $(1 - \epsilon)d/2$ clauses w.r.t. to MAJ (the clauses which prevented the removal of $x$ from $H$ in step 2 in the construction of $H$). Hence $H$ survives the first round of unassignment (lines 3-5). This invariant is kept throughout the unassignment step to prove that $H$ remains assigned at its end. Since this (lines 3-5) is the only place where unassignment is performed, the lemma follows. ■

**Proof.** (Proposition 11) Let $K$ be a **rv** measuring the size of the largest connected component in $G(\bar{H}, E)$. Lemma 12 implies that $U \subseteq \bar{H}$. Therefore, $R_{EXH} \leq n \cdot 2^K \cdot O(K^3)$; $n$ upper bounds the number of connected components, $2^K$ is the number of possible assignments $\tau(\Gamma_i)$, and $O(K^3)$ is the time it takes to check for every $\tau(\Gamma_i)$ whether it satisfies $I|_{\psi \circ \pi(Y)}[\Gamma_i]$ (the number of clauses in $I|_{\psi \circ \pi(Y)}[\Gamma_i]$ is $O(K^3)$). We introduce some notations to simplify the presentation.

$$f(k, \alpha) = n \cdot 2^k \cdot O(k^3) \cdot \binom{n}{\alpha n} \cdot 2^{\alpha n}.$$

$$g(k, \alpha) = Pr[K = k, |Y_0| = \alpha n].$$

$$E[R_{REC} \cdot R_{EXH}] \leq \sum_{k,\alpha} f(k, \alpha) \cdot g(k, \alpha).$$

Let $k_0 = \Theta(\log n/(\epsilon^6 d))$. For $k \geq k_0, \alpha \leq \alpha_0$,

$$g(k, \alpha) \leq \exp\{-\Omega(\max\{\underbrace{d \cdot k}_{\text{Proposition 6}}, \underbrace{d \cdot \log \frac{1}{\alpha} \cdot \alpha n}_{\text{Prop. 4}}\})\}.$$

For $k \geq k_0, \alpha \geq \alpha_0$,

$$g(k, \alpha) \leq \exp\{-\max\{\Omega(d \cdot k), \underbrace{d/4000 \cdot \alpha n}_{\text{Prop. 2}}\}\}$$

and

$$f(k, \alpha) \leq \exp\{\Theta(k + d/5000 \cdot \alpha n)\}.$$

Similarly to the proof of Proposition 9, the proposition follows by breaking the sum according to the different values of $\alpha$ and $k$:

$$\sum_{k<k_0}\sum_{\alpha} f(k,\alpha) \cdot g(k,\alpha) \le n \cdot 2^{k_0} \cdot O(k_0^3) \cdot \sum_{\alpha} Pr[|Y_0| = \alpha n] \cdot \exp\{(\log \frac{1}{\alpha} + 3) \cdot \alpha n\}$$

$$\underbrace{=}_{\text{Prop. 9}} n^{1+\Theta(1/d)} \cdot O(1) \ (a)$$

$$\sum_{k_0 \le k}\sum_{\alpha:k<\log \frac{1}{\alpha} \cdot \alpha n} f(k,\alpha) \cdot g(k,\alpha) \le \sum_{k_0 \le k}\sum_{\alpha} \exp\{-\Omega(d) \cdot \alpha n\} \le \sum_{k_0 \le k} 2 \le 2n \ (b)$$

$$\sum_{k_0 \le k}\sum_{\alpha:\log \frac{1}{\alpha} \cdot \alpha n \le k} f(k,\alpha) \cdot g(k,\alpha) \le \sum_{k_0 \le k}\sum_{\alpha} \exp\{-\Omega(d) \cdot k\} = \sum_{k} \exp\{-\Omega(d) \cdot k\} \sum_{\alpha} 1 \le 2n \ (c)$$

Since $E[R_{REC} \cdot R_{EXH}] \le (a) + (b) + (c)$, the proposition follows. $\blacksquare$

Combining Propositions 7, 8, 9, and 11 one obtains $E[R(\mathcal{A})] \le n^{1+\Theta(1/d)}$ as desired. This also finishes the proof of Theorem 1.

# 5 Discussion

To summarize, we presented an algorithm that solves 3SAT. Of course, one cannot expect the algorithm to be polynomial, however, over the planted distribution with $p \ge d/n^2$, the algorithm has an expected polynomial running time. Previous work on planted 3SAT present algorithms that work **whp**. Our result extends the latter by presenting an expected polynomial time algorithm. An expected polynomial time algorithm also works **whp**, but in addition it works for all instances, while keeping an overall average polynomial running time. A few interesting issues are open for future research.

## 5.1 The $k$-opt heuristic

In [FV04], a different algorithm, based on the $k$-opt heuristic, is suggested. A succinct description follows. The algorithm can be viewed as performing a search in the graph whose vertices are all $2^n$ assignments, and two assignments are connected if they differ by the assignment of exactly one variable. The algorithm starts at the vertex corresponding to the majority vote on $I$. The goal is to reach a vertex corresponding to a satisfying assignment. We say that an assignment $\psi'$ improves $\psi$ if the set of clauses satisfied by $\psi'$ contains the set of clauses satisfied by $\psi$. While not at a satisfying assignment, in every iteration the algorithm checks all possible improving assignments at distance at most $k$. If such is found, the algorithm moves to that vertex. Otherwise, the algorithm has reached a local optimum, and the algorithm fails (or in a different version, $k$ is increased by 1 and the algorithm continues with the new $k$). It is shown in [FV04] that an improving assignment at distance at most $k$, if such exists, can be found in time $O(n2^k)$. Thus, by demanding $k = O(\log n)$, the algorithm is polynomial. As the analysis in [FV04] shows, for $\mathbb{P}_{d/n^2,n}$, the algorithm finds a satisfying assignment **whp** when $k = \log n$.

An interesting question is whether one can transform the $k$-opt version to run in expected polynomial time. The difficulty lies in the fact that no explicit distinction between correctly and wrongly assigned variables is done.

## 5.2 Simplifying the coloring algorithm for $k$-colorable graphs

Let us compare our algorithm with the one presented in [Böt05] for coloring planted $k$-colorable graphs in expected polynomial time. The coloring algorithm is arguably more complicated than the satisfiability one. However, this is not the case when comparing [AK97] to [Fla03]. Why is the satisfiability algorithm simpler? We derive the first approximation to the planted solution by the majority vote. The coloring algorithm approximates the planted solution using SDP. Analogously to the set $H$, one can define such a set for planted $k$-colorable graphs. The role of the set $H$ is crucial in the analysis of both algorithms. One first proves that **whp** the algorithm sets $H$ according to the planted solution, allowing $H$ to survive the unassignment step. Further, one proves that **whp** $H$ is sufficiently large so that the residual graph $G(\bar{H}, E)$ has no connected component of size larger than $\log n$. Therefore, an exhaustive search can be employed to complete the solution. The analysis in [Böt05] doesn't show that the SDP approximation sets $H$ correctly with sufficiently high probability to allow an expected polynomial time algorithm. Hence, another refinement step precedes the unassignment. In our case, $H$ is set correctly by definition, and the refinement step can be skipped. This simplifies greatly the algorithm and the analysis. Further, the majority vote is much simpler to analyze and implement compared to SDP. An interesting research question is replace the SDP step in the coloring algorithm with a simpler and more natural step analogous to the majority vote.

# References

[AK97]   Noga Alon and Nabil Kahale. A spectral technique for coloring random 3-colorable graphs. *SIAM J. on Comput.*, 26(6):1733–1748, 1997.

[AKS98]  Noga Alon, Michael Krivelevich, and Benny Sudakov. Finding a large hidden clique in a random graph. *Random Structures and Algorithms*, 13(3-4):457–466, 1998.

[AS00]   Noga Alon and Joel H. Spencer. *The probabilistic method.* Wiley-Interscience Series in Discrete Mathematics and Optimization. Wiley-Interscience [John Wiley & Sons], New York, second edition, 2000. With an appendix on the life and work of Paul Erdős.

[AZ04]   Martin Aigner and Günter M. Ziegler. *Proofs from The Book.* Springer-Verlag, Berlin, third edition, 2004. Including illustrations by Karl H. Hofmann.

[Böt05]  Julia Böttcher. Coloring sparse random $k$-colorable graphs in polynomial expected time. In *Proc. 30th International Symp. on Mathematical Foundations of Computer Science*, volume 2461 of *Lecture Notes in Comput. Sci.*, pages 156–167, 2005.

[BS95]   Avrim Blum and Joel Spencer. Coloring random and semi-random $k$-colorable graphs. *J. of Algorithms*, 19(2):204–234, 1995.

[BSBG02] Eli Ben Sasson, Yonatan Bilu, and Dan Gutfreund. Finding a randomly planted assigment in a random $3CNF$. *manuscript*, 2002.

[BV04]      Rene Beier and Berthold Vöcking. Random knapsack in expected polynomial time. *J. Comput. and Syst. Sci.*, 69(3):306–329, 2004.

[BW85]      Edward A. Bender and Herbert S. Wilf. A theoretical analysis of backtracking in the graph coloring problem. *J. of Algorithms*, 6(2):275–282, 1985.

[CA96]      James M. Crawford and Larry D. Auton. Experimental results on the crossover point in random 3-SAT. *Artificial Intelligence*, 81(1-2):31–57, 1996.

[CI95]       Byungki Cha and Kazuo Iwama. Performance test of local search algorithms using new types of random CNF formulas. In *Proc. 14th International Joint Conference on Artificial Intelligence*, pages 304–311, 1995.

[CO03]      Amin Coja-Oghlan. Finding large independent sets in polynomial expected time. In *Proc. 20th Symp. on Theoretical Aspects of Comp. Science*, volume 2607 of *Lecture Notes in Comput. Sci.*, pages 511–522. Springer, Berlin, 2003.

[CO04]      Amin Coja-Oghlan. Coloring semirandom graphs optimally. In *Proc. 31st International Colloquium on Automata, Languages, and Programming*, pages 383–395, 2004.

[COGLS04] Amin Coja-Oghlan, Andreas Goerdt, André Lanka, and Frank Schädlich. Techniques from combinatorial approximation algorithms yield efficient algorithms for random $2k$-SAT. *Theoret. Computer Sci.*, 329(1-3):1–45, 2004.

[Coo71]     Stephen A. Cook. The complexity of theorem-proving procedures. In *Proc. 3rd ACM Symp. on Theory of Computing*, pages 151–158, 1971.

[COT03]     Amin Coja-Oghlan and Anusch Taraz. Colouring random graphs in expected polynomial time. In *Proc. 20th Symp. on Theoretical Aspects of Comp. Science*, volume 2607 of *Lecture Notes in Comput. Sci.*, pages 487–498. Springer, Berlin, 2003.

[DF89]      Martin Dyer and Alan Frieze. The solution of some random NP-hard problems in polynomial expected time. *J. of Algorithms*, 10(4):451–489, 1989.

[FK00]      Uriel Feige and Robert Krauthgamer. Finding and certifying a large hidden clique in a semirandom graph. *Random Structures and Algorithms*, 16(2):195–208, 2000.

[Fla03]      Abraham Flaxman. A spectral technique for random satisfiable 3CNF formulas. In *Proc. 14th ACM-SIAM Symp. on Discrete Algorithms*, pages 357–363, 2003.

[FM97]      Alan Frieze and Colin McDiarmid. Algorithmic theory of random graphs. *Random Structures and Algorithms*, 10(1-2):5–42, 1997.

[Fri99]       Ehud Friedgut. Sharp thresholds of graph properties, and the $k$-sat problem. *J. Amer. Math. Soc.*, 12(4):1017–1054, 1999.

[FSVM93]  Martin Fürer, C. R. Subramanian, and C. E. Veni Madhavan. Coloring random graphs in polynomial expected time. In *Proc. 4th Annual International Symposium on Algorithms and Computation*, volume 762 of *Lecture Notes in Comput. Sci.*, pages 31–37. Springer, Berlin, 1993.

[FV04]     Uriel Feige and Dan Vilenchik. A local search algorithm for 3SAT. Technical report, The Weizmann Institute of Science, 2004.

[Gen98]    Ian Gent. On the stupid algorithm for satisfiability. Technical report, Strathclyde University, 1998.

[KKL02]    Alexis C. Kaporis, Lefteris M. Kirousis, and Efthimios G. Lalas. The probabilistic analysis of a greedy satisfiability algorithm. In *Proc. 10th Annual European Symposium on Algorithms*, volume 2461 of *Lecture Notes in Comput. Sci.*, pages 574–585. Springer, Berlin, 2002.

[KKS⁺01]  Alexis C. Kaporis, Lefteris M. Kirousis, Yannis C. Stamatiou, Malvina Vamvakari, and Michele Zito. Coupon collectors, $q$-binomial coefficients and the unsatisfiability threshold. In *Proc. Theoretical Computer Science: 7th Italian Conference*, volume 2202 of *Lecture Notes in Comput. Sci.*, pages 328–338. Springer, Berlin, 2001.

[KP92]     Elias Koutsoupias and Christos H. Papadimitriou. On the greedy algorithm for satisfiability. *Info. Process. Letters*, 43(1):53–55, 1992.

[Kri02]    Michael Krivelevich. Deciding $k$-colorability in expected polynomial time. *Info. Process. Letters*, 81(1):1–6, 2002.

[KV02]     Michael Krivelevich and Van H. Vu. Approximating the independence number and the chromatic number in expected polynomial time. *J. Comb. Optim.*, 6(2):143–155, 2002.

[PS92]     Hans Jürgen Prömel and Angelika Steger. Coloring clique-free graphs in linear expected time. *Random Structures and Algorithms*, 3(4):375–402, 1992.

[Vil04]    Dan Vilenchik. Finding a satisfying assignment for semi-random satisfiable $3CNF$ formulas. Master's thesis, The Weizmann Institute of Science, Rehovot, Israel, January 2004.

[Wil84]    Herbert S. Wilf. Backtrack: an $O(1)$ expected time algorithm for the graph coloring problem. *Info. Process. Letters*, 18(3):119–121, 1984.

# A    Proof of Proposition 2

For a variable $x$, let $F_x$ be the event $\{MAJ(x) \neq \varphi(x)\}$. Assume w.l.o.g that $\varphi(x) = T$. Let $P_x$ (resp. $N_x$) be a **rv** counting the number of clauses in which $x$ appears positively (resp. negatively). It is readily seen that

$$P_x \sim Binom\left(7\binom{n-1}{2}, 4/7 \cdot d/n^2\right)$$

$$N_x \sim Binom\left(7\binom{n-1}{2}, 3/7 \cdot d/n^2\right).$$

Therefore,

$$E[P_x] = 2d + o(1), E[N_x] = 3d/2 + o(1).$$

Namely there is a gap of $d/2$ in the expectations. If the majority is wrong for $x$ then at least one of $N_x, P_x$ has deviated from its expectation by at least $d/4$. Following the intuition given in Proposition 2, our first goal is to bound $Pr[F_{x_1} \wedge F_{x_2} \wedge ... \wedge F_{x_{\alpha n}}]$ for $\alpha \in [\alpha_0, \epsilon]$.

**Proposition 13.** *Let $T = \{x_1, x_2, .., x_{\alpha n}\}$ be a fixed set of variables. $Pr[F_{x_1} \wedge F_{x_2} \wedge ... \wedge F_{x_{\alpha n}}] \leq exp\{-d/2000 \cdot \alpha n\}$*

Before proving this proposition, some more observations are due. Let $f(T)$ be the number of clauses in $I$ containing at least two variables from $T$. It is easily seen that

$$E[f(T)] \leq 7 \cdot \binom{|T|}{2} \cdot n \cdot d/n^2 \leq 7\alpha d/2 \cdot \alpha n.$$

**Lemma 14.** *For $T$ as above, $Pr[f(T) \geq d/60 \cdot \alpha n] \leq \exp\{-d/80 \cdot \alpha n\}$.*

The proof of this Lemma consists of standard probabilistic arguments (Chernoff bound for example). Assume indeed that $f(T) \leq d/60 \cdot \alpha n$. We say that a variable in $T$ is *T-influenced* if it appears in more that $d/10$ clauses with at least another variable from $T$. Let $Y \subseteq T$ be the set of $T$-influenced variables. It holds that $(|Y| \cdot d/10)/3 \leq f(T)$ (since every clause was counted at most 3 times). Put differently, $|Y| \leq 30f(T)/d \leq \alpha n/2$ (by our assumption on $f(T)$). Now we can approach the proof of Proposition 13.

**Proof.** (Proposition 13) Assume that $f(T) \leq d/60 \cdot \alpha n$. W.l.o.g, let $x_1, .., x_{\alpha n/2}$ be a set of non-$T$-influenced variables (namely a subset of $T \setminus Y$). Let $\hat{P}_{x_i}$ (resp. $\hat{N}_{x_i}$) be a **rv** counting the positive (resp. negative) appearances of $x_i$ when only clauses where $x_i$ is the only variable from $T$ are counted. Assuming $\varphi(x_i) = T$,

$$E[\hat{P}_{x_i}] = 4 \cdot \binom{(1-\alpha)n}{2} \cdot d/n^2 = (1-\alpha)^2 2d + o(1)$$

$$E[\hat{N}_{x_i}] = (1-\alpha)^2 3d/2 + o(1).$$

Since $\alpha \leq \epsilon$, it follows then that

$$|E[\hat{P}_{x_1}] - E[\hat{N}_{x_1}]| > d/3.$$

For $i \in [1, \alpha n/2]$, if $F_{x_i}$ occurs, then at least one of $\hat{P}_{x_i}, \hat{N}_{x_i}$ has deviated from its expectation by more than $(d/3 - d/10)/2 > d/10$ (otherwise the majority is not wrong, since $x_i$ appears in at most $d/10$ "misleading" clauses with variables from $T$). Using the Chernoff bound once again,

$$Pr[F_{x_i}] \leq Pr[|\hat{P}_{x_i} - E[\hat{P}_{x_i}]| \geq d/10] + Pr[|\hat{N}_{x_i} - E[\hat{N}_{x_i}]| \geq d/10] \leq \exp\{-d/900\}.$$

Observe that all the $\hat{P}_{x_i}, \hat{N}_{x_i}, i \in [1, \alpha n]$ are mutually independent, as they correspond to disjoint sets of clauses (by their definition). For brevity, Let $F^j = F_{x_1} \wedge F_{x_2} \wedge ... \wedge F_{x_j}$. To conclude,

$$Pr[F^{\alpha n}] \leq Pr[F^{\alpha n/2}|f(T) \leq d/60 \cdot \alpha n] \cdot Pr[f(T) \leq d/60 \cdot \alpha n] + Pr[f(T) \geq d/60 \cdot \alpha n]$$

$$\leq \exp\{-d/900 \cdot \alpha n/2\} \cdot (1 - \exp\{-d/80 \cdot \alpha n\}) + \exp\{-d/80 \cdot \alpha n\} \leq \exp\{-d/2000 \cdot \alpha n\}$$

$\blacksquare$

To conclude the proof of Proposition 2, taking the union bound over all possible sets $T$,

$$Pr[F_{MAJ} \geq \alpha n] \leq \binom{n}{\alpha n} \cdot \exp\{-d/2000 \cdot \alpha n\} \leq$$

$$\leq (e/\alpha)^{\alpha n} \cdot \exp\{-d/2000 \cdot \alpha n\} \leq \exp\{\alpha n(1 + \log \frac{1}{\alpha} - d/2000) \leq \exp\{-d/4000 \cdot \alpha n\}$$

The last inequality is due to $\alpha \geq \alpha_0 = \exp\{-d/5000\}$.

# B   proof of Proposition 6

**Proof.** Our strategy is to show that **whp**, the largest tree in $G(\bar{H}, E)$ is of size $\beta n$. Let $T$ be a fixed tree on $\beta n$ variables, and let $T'$ be a fixed collection of clauses such that each edge of $T$ is induced by some clause of $T'$. Let $V(T)$ be the set of vertices of $T$, and $E(T)$ be the set of edges of $T$. We say that an edge of $T$ is *covered* by $T'$ if there exists a clause in $T'$ that contains both its endpoints. We say that $T'$ is *minimal* if by deleting a clause from $T'$, $T$ is not covered by $T'$ anymore. By the definition of minimality, $|T'| \leq |E(T)| = |V(T)| - 1$ ($T$ is a tree). Our first goal is to bound

$$Pr[T' \subseteq I \text{ and } V(T) \cap H = \emptyset].$$

Both $Pr[T' \subseteq I]$ and $Pr[V(T) \cap H = \emptyset]$ are easily calculated (for fixed $T, T'$). However, since $H$ is not a truly random subset of $V$, the events $\{T' \subseteq I\}$ and $\{V(T) \cap H = \emptyset\}$ are not independent. To decouple the dependency, we introduce two new sets, $J$ and $H'$, replacing $V(T)$ and $H$.

The analysis follows closely the one given in [Fla03], [Vil04]. However, since we describe an expected polynomial time algorithm, a more tight analysis is due. We omit most of the proofs, given fully in [Fla03], [Vil04] and only point out the differences.

Let $J \subseteq V(T)$ be a set of variables which appear in at most 6 clauses of $T'$. For a set of clauses $I$, let $|I|$ denote the number of clauses in $I$.

**Proposition 15.** $|J| \geq \frac{|V(T)|}{2}$

**Proof.** Suppose that $|J| < \frac{|V(T)|}{2}$, then more than $\frac{|V(T)|}{2}$ variables of $T'$ appear in more than 6 clauses. All in all they appear in at least $6 \cdot \frac{|V(T)|}{2}$ clauses. How many clauses were counted more than once? every clause was counted at most 3 times, once for every variable appearing in it. Hence, there are at least $\frac{1}{3} \cdot 3 \cdot |V(T)|$ distinct clauses in $T' \Rightarrow |T'| \geq |V(T)|$.
On the other hand, $|T'| \leq |V(T)| - 1$, contradicting the above. ∎

We define the set of variables $H'$ by the following iterative procedure:
*Let $V$ denote the set of all variables.*
*Let $W$ be the set of variables for which the majority vote is wrong or for which the majority is right with advantage smaller than 7.*
*Let $B$ be the set of variables which support less than $(1 - \epsilon/2)d/2$ clauses w.r.t. $\varphi$.*
*Let $H' \subset V$ be constructed as follows:*

1. *Let $H'_0 = V \setminus \{W \cup B \cup (V(T') \setminus J)\}$.*

2. *While there is a variable $a_i \in H'_i$ that supports less than $\frac{d}{2} - \epsilon d$ clauses with only variables of $H'_i$, define $H'_{i+1} = H'_i \setminus \{a_i\}$.*

3. *Let $a_m$ be the last variable removed at step 2. Define $H' = H'_{m+1}$.*

The difference between the procedure for $H$ and $H'$ is in the definition of $W$ and in the definition of $H'_0$.

**Claim 16.** *Let $\overline{H'} = V \setminus H'$, let $\beta \leq \alpha \in [4\alpha_1, (\epsilon/128)^4]$. $Pr[\overline{H'} \geq 4\alpha n] \leq \exp\{-\epsilon^2 d/100 \cdot \alpha n\}$.*

**Proof.** The proof is essentially the same as Proposition 5. However, two things need extra care, the size of $V(T') \setminus J$ and the size of $W$. We know that $|V(T') \setminus J| \leq |V(T')| \leq 3|V(T)| \leq 3\beta n \leq 3\alpha n$. The result of Proposition 2 can be repeated when demanding a gap of 7 in the majority vote, since enough slackness was left for that in the choice of constants in Proposition 2. Now if $\overline{H}' \geq 4\alpha n$ then at least one of $\{W \cup B\}$ or $\overline{H}^{it}$ is of size at least $(4\alpha n - 3\alpha n)/2 = \alpha n/2$ and the proof of Proposition 5 is repeated (again, enough slackness in the choice of constants in the bound was left to accommodate $\overline{H}'$ rather than $\overline{H}$). ∎

**Claim 17.** $Pr[T' \subset I \text{ and } V(T) \cap H = \emptyset] \leq Pr[T' \subset I]Pr[J \cap H' = \emptyset]$.

The key for proving the proposition is in the following inclusion.

**Lemma 18.** *Let $F$ be a set of clauses and let $H(F \cup T')$ be the value of $H$ if $I = F \cup T'$. Let $H'(F)$ be the value of $H'$ if $I = F$. Then $H'(F) \subseteq H(F \cup T')$.*

This lemma clarifies the motivation for defining $H'$. It is not true that $H(F) \subseteq H(F \cup T')$ for all $F$. This is because a variable which appeared in $H(F)$ could disappear from $H(F \cup T')$ since the clauses of $T'$ would make the majority vote wrong for it (however, they cannot reduce its support). We overcome this problem by deleting from $H'$ vertices with high degree in $T'$ (in our case, vertices with degree greater than 6). This on the one hand reduces the effect of $T'$ on $H'$, but on the other leaves enough vertices in $J$ (at least $\frac{|V(T)|}{2}$) such that the value $Pr[J \cap H' = \emptyset]$ would be sufficiently small to help bound the probability for a tree.

**Proof.**(Lemma 18) The lemma is proved using induction. First we show that $H'_0(F) \subseteq H_0(F \cup T')$. If $x \in H'_0(F)$ then the majority vote for $x$ is right with advantage of at least 7, and $x$ supports at least $\frac{d}{2} - \epsilon d$ clauses. Since the degree of $x$ in $T'$ is at most 6, the majority vote could not be flipped by adding $T'$ to $F$. Also the support could only increase by adding $T'$, hence $x \in H_0(F \cup T')$ as well. We proceed by showing that during the iterative step we still have $H'_i(F) \subseteq H_i(F \cup T)$ for every $i$.

Suppose that $H'_i(F) \subseteq H_i(F \cup T')$. We show that the lemma holds for $i+1$. If $x$ supports at least $\frac{d}{2} - 2\epsilon d$ clauses with only variables of $H'_i(F)$ w.r.t. $F$, then since $H'_i(F) \subseteq H_i(F \cup T')$ and $F \subseteq F \cup T'$, it supports the same clauses w.r.t. $(F \cup T')$ with only variables of $H_i(F \cup T')$. Hence if it remains in $H'_{i+1}(F)$ it will also remain in $H_{i+1}(F \cup T')$. ∎

**Proof.**(Proposition 17)

$$Pr[T' \subset I \text{ and } V(T) \cap H = \emptyset] \underbrace{\leq}_{J \subseteq V(T)} Pr[T' \subset I \text{ and } J \cap H = \emptyset] = Pr[J \cap H = \emptyset | T' \subset I]Pr[T' \subset I].$$

Therefore, it suffices to prove

$$Pr[J \cap H = \emptyset | T' \subset I] \leq Pr[J \cap H' = \emptyset].$$

$$Pr[J \cap H' = \emptyset] = \sum_{F : J \cap H'(F) = \emptyset} Pr[I = F] \underbrace{\geq}_{H'(F) \subseteq H(F \cup T')} \sum_{F : J \cap H(F \cup T') = \emptyset} Pr[I = F]$$

Break each set of clauses $F$ into $F' = F \setminus T'$ and $F'' = F \cap T'$, and the latter equals

$$\sum_{F' : F' \cap T' = \emptyset, J \cap H(F' \cup T') = \emptyset} \sum_{F'' : F'' \subseteq T'} Pr[I \setminus T' = F' \wedge I \cap T' = F'']$$

18

Since the two sets of clauses, $I \setminus T'$, and $I \cap T'$, are disjoint, and clauses are chosen independently, the last expression equals,

$$\sum_{F':F'\cap T'=\emptyset, J\cap H(F'\cup T')=\emptyset} \sum_{F'':F''\subseteq T'} Pr[I\setminus T'=F']Pr[I\cap T'=F''] =$$

$$\sum_{F':F'\cap T'=\emptyset, J\cap H(F'\cup T')=\emptyset} Pr[I\setminus T'=F'] \underbrace{\sum_{F'':F''\subseteq T'} Pr[I\cap T'=F'']}_{1} =$$

$$\sum_{F':F'\cap T'=\emptyset, J\cap H(F'\cup T')=\emptyset} Pr[I\setminus T'=F']$$

Since $(I\setminus T')\cap T' = \emptyset$, and clauses are chosen independently, the event $\{T'\subseteq I\}$ is independent of the event $\{I\setminus T'=F'\}$. Therefore, the latter expression can be rewritten as

$$\sum_{F':F'\cap T'=\emptyset, J\cap H(F'\cup T')=\emptyset} Pr[I\setminus T'=F'|T'\subseteq I] = Pr[J\cap H=\emptyset|T'\subseteq I].$$

■

**Lemma 19.** *Let $\rho \in [0,1]$. Then, $\binom{\rho n}{k}/\binom{n}{k} \le (6\rho)^k$*

**Proof.** $\binom{\rho n}{k} \le (e\rho n/k)^k$, $\binom{n}{k} \ge (n-k/k)^k \ge (n/2k)^k$. Hence, $\binom{\rho n}{k}/\binom{n}{k} \le (e\rho n/k)^k/(n/2k)^k \le (6\rho)^k$

■

**Lemma 20.** *Let $\rho \in [\beta, 1]$ be s.t. $|H'| = (1-\rho)n$. Then, $Pr[J\cap H'=\emptyset] \le (6\rho)^{\beta n/2}$*

**Proof.** Since the choice of $H'$ is independent of $J$, it follows that

$$Pr[J\cap H'=\emptyset] \le \frac{\binom{n-|J|}{|H'|}}{\binom{n}{|H'|}} = \frac{\binom{n-|H'|}{|J|}}{\binom{n}{|J|}} = \frac{\binom{\rho n}{\beta n/2}}{\binom{n}{\beta n/2}} \le (6\rho)^{\beta n/2}.$$

The last inequality follows from the previous Lemma. ■

**Corollary 21.** *Let $\rho$ be as in the previous Lemma. Then, $Pr[T' \subset I \wedge V(T)\cap H = \emptyset] \le (6\rho)^{|T'|/2}(d/n^2)^{|T'|}$*

**Proof.** By Proposition 17, $Pr[T' \subset I \wedge V(T)\cap H = \emptyset] \le Pr[T' \subset I]Pr[J\cap H' = \emptyset]$. By the previous lemma, $Pr[J\cap H' = \emptyset] \le (6\rho)^{|T'|/2}$. Since the clauses of $I$ are chosen independently, $Pr[T' \subset I] = (d/n^2)^{|T'|}$. ■

**Lemma 22.** *Let $T$ be a fixed tree with $k$ vertices. The number of minimal sets of clauses that span $T$ is at most*

$$\sum_{s=0}^{k/2} \binom{k}{2} 7^{k-s-1} n^{k-2s-1}$$

19

**Proof.** Let $T'$ be a minimal spanning set. We proved that a minimal set contains at most $k-1$ clauses. Every clause may cover up to 2 edges (otherwise $T$ contains a cycle). We can extend a tree to a minimal set by having $s$ clauses cover 2 edges of $T'$ and $k-1-2s$ clauses cover one edge. The total number of clauses in $T'$ is $k-s-1$. For every clause there are $(2^3-1)$ ways to set the signs of the variables. For every clause that covers one edge, there are $(n-2)$ ways to choose the third variable. A clause that covers 2 edges is determined uniquely by the two edges. We have at most $\binom{k}{2}$ ways to form the clauses that cover 2 edges. $s$ may range from 0 to $k/2$. ∎

**Theorem 23.** *(Cayley) The number of spanning trees of $K_n$ is $n^{n-2}$*

Four different proofs are given at [AZ04].

**Claim 24.** *The probability of a tree of size $\beta n$, $\beta n = \Omega(\log n / n\epsilon^6 d)$, in $G(\overline{H}, E)$ is at most $\exp\{-\Omega(\epsilon^6 d)\cdot\beta n\}$*

**Proof.** Set $\alpha_2 = 4(\epsilon/128)^4$.

$Pr[\text{a tree of size } \beta n \text{ in } G(\overline{H}, E)] \leq Pr[\exists T' \subset I, |V(T)| = \beta n \wedge V(T)\cap H = \emptyset] =$

$$\sum_{\rho=\beta}^{1} Pr[\exists T' \subset I, |V(T)| = \beta n \wedge V(T)\cap H = \emptyset \mid |\overline{H'}| = \rho n]\cdot Pr[|\overline{H'}| = \rho n] \underbrace{\leq}_{\text{union bound}}$$

$$\sum_{\rho=\beta}^{1} Pr[|\overline{H'}| = \rho n] \sum_{s=0}^{\beta n/2} \underbrace{\binom{n}{\beta n}\cdot(\beta n)^{\beta n-2}}_{\text{number of trees of size } \beta n} \cdot \underbrace{7^{\beta n-s-1}\cdot n^{\beta n-2s-1}\binom{\beta n}{2}}_{\text{extend to a minimal set}} \underbrace{\cdot(d/n^2)^{\beta n-s-1}\cdot(6\rho)^{\beta n/2}}_{\text{variables don't intersect } H \text{ (cor. 21)}} \leq$$

$$\sum_{\rho=\beta}^{\max(16\alpha_1,\beta)} (6\rho)^{\beta n/2}(e\cdot d)^{\beta n}\beta n/2 + \sum_{\rho=\max(16\alpha_1,\beta)}^{1} Pr[|\overline{H'}| = \rho n] \underbrace{=}_{\text{Prop. 16}}$$

$$\sum_{\rho=\beta}^{\max(16\alpha_1,\beta)} (6\rho)^{\beta n/2}(e\cdot d)^{\beta n}\beta n/2 + \sum_{\rho=\max(16\alpha_1,\beta)}^{\max(\alpha_2,\beta)} \exp\{-\epsilon^2 d/100\cdot\rho n\} + \sum_{\rho=\max(\alpha_2,\beta)}^{1} \exp\{-\Omega(\epsilon^6 d)n\}$$

$$\leq \sum_{\rho=\beta}^{\max(16\alpha_1,\beta)} \exp\{-\Omega(\epsilon^2 d)\cdot\beta n\} + \sum_{\rho=\max(16\alpha_1,\beta)}^{\max(\alpha_2,\beta)} \exp\{-\epsilon^2 d/100\cdot\beta n\} + \sum_{\rho=\max(\alpha_2,\beta)}^{1} \exp\{-\Omega(\epsilon^6 d)\beta n\}$$

$$\leq n\cdot\exp\{-\Omega(\epsilon^6 d)\cdot\beta n\} = \exp\{-\Omega(\epsilon^6 d)\cdot\beta n + \log n\} \underbrace{=}_{\text{choice of }\beta} \exp\{-\Omega(\epsilon^6 d)\cdot\beta n\}$$

∎

Proposition 6 follows easily from the last proposition. ∎

# C  Proof of Proposition 8

Let $m$ be a random variable counting the number of clauses in $I$. As calculated above, $E[m] = d/n^2\cdot 7\cdot\binom{n}{3} = 7dn/6 + o(n) = O(dn)$. To prove Proposition 8, we present procedures that calculate

the majority vote, perform the unassignment, compute the connected components of $G(U, E)$ and simplify $I$ in time $O(m + n)$. Since $E[O(m + n)] = O(E[m + n]) = O(dn)$, the proposition follows.

To calculate the majority vote in time $O(m + n)$, allocate an array $A$ of size $n$, and reset all entries. Go over the instance $I$ once, for every clause $C = (l_i \vee l_j \vee l_k)$, if $l_i = x_i$ perform $A[i] \leftarrow A[i] + 1$, otherwise $(l_i = \overline{x_i})$ $A[i] \leftarrow A[i] - 1$. Same applies to $l_j, l_k$. Next, go over $A$, if $A[i] \geq 0$ then assign $x_i$ true, otherwise false.

To perform the unassignment step in time $O(m + n)$, we assume that $I$ is represented as an adjacency list. Namely, an array $A[1..n]$, s.t. $A[i]$ is a list of all the clauses containing $x_i$. We keep another array $S[1..n]$, and a list $U$, which is empty at the beginning. In the first step, calculate for every $x_i$ its support and store it in $S[i]$. If the support of $x_i$ is less than $(1 - \epsilon)d/2$, add $x_i$ into $U$. The first step consumes time $O(m + n)$ (as each clause is scanned at most 3 times). Next, as long as $U \neq \emptyset$, take out a variable $x_k \in U$ and unassign it; go over the list of clauses in $A[k]$, for every clauses $C$, if there exists a variable $x_j$ that supported $C$ (w.r.t. to the current partial assignment), update $S[j] \leftarrow S[j] - 1$. If $S[j] < (1 - \epsilon)d/2$, add $x_j$ to $U$. The second step consumes time $O(m)$ as each clause is scanned at most 3 times.

Computing the connected components of $G(U, E)$ can be done in time $O(|U| + |E|) \leq O(m + n)$ using BFS.

As for the simplification procedure, for a partial assignment $\psi$, one can calculate $I|_\psi$ by going over all clauses once and simplifying each clause. This takes time $O(m)$. One can calculate $I|_\psi[\Gamma_i]$ (the subformula of $I|_\psi$ induced on $\Gamma_i$) for all $\Gamma_i$'s by going over $I|_\psi$ once, and decide for every clause to which $I|_\psi[\Gamma_i]$ it belongs (in constant time). All in all, the running time is $O(m)$.