

ELIMINATION OF INFREQUENT VARIABLES IMPROVES AVERAGE CASE PERFORMANCE OF SATISFIABILITY ALGORITHMS*

JOHN FRANCO[†]

Abstract. We consider pre-processing a random instance I of CNF Satisfiability in order to remove infrequent variables (those which appear once or twice in an instance) from I . The model used to generate random instances is the popular random-clause-size model with parameters n , the number of clauses, r , the number of Boolean variables from which clauses are composed, and p , the probability that a variable appears in a clause as a positive (or negative) literal. It is shown that exhaustive search over such pre-processed instances runs in polynomial average time over a significantly larger parameter space than has been shown for any other algorithm under the random-clause-size model when $n = r^\epsilon$, $\epsilon < 1$, and $pr < \sqrt{\epsilon r \ln(r)}$. Specifically, the results are that random instances of Satisfiability are “easy” in the average case if $n = r^\epsilon$, $2/3 > \epsilon > 0$, and $pr < (\ln(n)/4)^{1/3} r^{2/3 - \epsilon}$; or $n = r^\epsilon$, $1 > \epsilon \geq 2/3$, $pr < (1 - \epsilon - \delta) \ln(n)/\epsilon$ for any $\delta > 0$; or $pn \rightarrow 0$, $pr < \gamma \ln \ln(n)$ for any $\gamma > 0$.

Key words. Satisfiability, NP-complete, Probabilistic Analysis, Resolution

1. Introduction. The Satisfiability problem is to determine whether there exists a truth assignment to the variables of a given CNF Boolean expression which cause it to have value *true*. If such a truth assignment exists we say the expression is satisfiable, otherwise it is unsatisfiable. The problem is NP-complete so there is no known polynomial time algorithm for solving it. Several papers have been concerned with the analysis of algorithms for Satisfiability that run in polynomial average time. These results depend on an assumed probabilistic input model. One popular model is the “random-clause-size” model which we refer to as $M(n, r, p)$.

Let $L = \{v_1, \bar{v}_1, v_2, \bar{v}_2, \dots, v_r, \bar{v}_r\}$ be a set of $2r$ literals. According to the model $M(n, r, p)$, n disjunctions (called clauses) are generated as follows: for each clause C_i , for all literals $l \in L$, put l in C_i with probability p , independently of the placement of other literals and clauses. Notice that it is possible for a pair of complementary literals (associated with the same variable) to be present in a clause. It is also possible to generate an empty (or null) clause using this model. If an instance contains a null clause it is trivially unsatisfiable. The preponderance or absence of null clauses in random instances is controlled by the product pr which is half the average number of literals in a clause. From [2] a random instance possesses a null clause with probability tending to 1 if the product $pr < \ln(n)/2$.

In the literature, polynomial average time results for Satisfiability algorithms are known only if $n = r^\epsilon$, $1 > \epsilon > 0$, $pr < \sqrt{\ln(n)} \cdot r^{1/2 - \epsilon}$ [4]; or $n = r^\gamma$, $\gamma > 1$, $pr < (\gamma - 1) \ln(n)/(2\gamma)$ [5]; or $n = \beta r$, β a positive constant, and $pr < f(\beta)$ where f is some complicated function of β ; or $pr > \sqrt{r \ln(n)}$ [3]. This space of parameters is depicted as region **I** in Figure 1 (the scale of the figure is such that factors of $\ln(n)$ are not distinguished).

* This work was carried out in part at the FAW, Helmholtzstraße 16, D-7900 Ulm/Donau, Germany. This work is based on research supported in part by the Air Force Office of Scientific Research, Grant No. AFOSR 84-0372 and 89-0186.

[†] Department of Computer Science, University of Cincinnati, Cincinnati, Ohio 45221-0008. john.franco@uc.edu or franco@ucunix.san.uc.edu

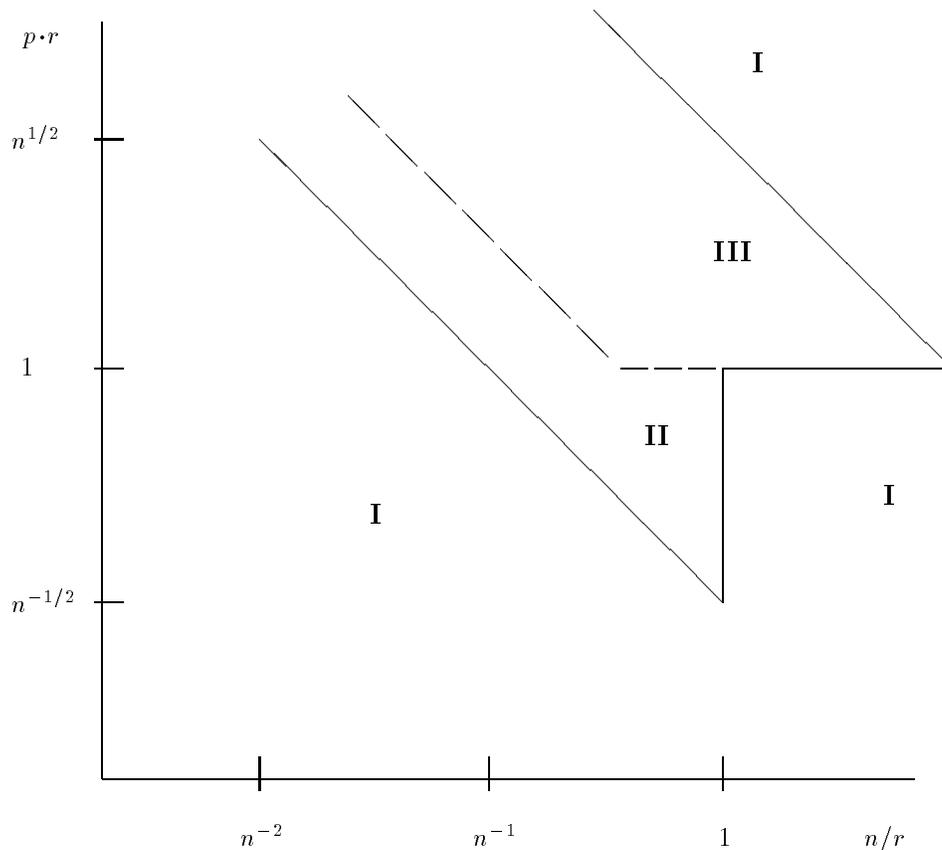


Fig 1. *Regions of the parameter space of model $M(n, r, p)$. Random instances generated using parameters which fall into region **I** are solved in polynomial time by some previously analyzed algorithm. The algorithm *PLR* requires superpolynomial average time for parameters set in regions **II** and **III**. Algorithm *INFREQ* solves random instances in polynomial average time in region **II**. No known algorithm takes polynomial average time for parameters set in region **III**.*

Furthermore, no polynomial average time results are published for parameters set in region **II** or region **III** of Figure 1. Also of interest is a result in [1] which shows that an algorithm based on the pure literal rule, called *PLR*, requires superpolynomial average time if $n = r^\epsilon$, $1 > \epsilon > 1/2$, and $pr > \sqrt{\omega(r) \ln(n)} \cdot r^{1/2-\epsilon}$ where $\omega(r)$ is any growing function of r . In Figure 1, this range of parameters includes region **II**. Thus, there is a significant range of pr for which null clauses exist in random instances with high probability but no published polynomial average time analysis exists and at least one non-trivial algorithm, namely *PLR*, requires superpolynomial average time. This range is depicted in Figure 1 as the part of region **II** below the extension of the lower boundary between regions **I** and **III**.

In this paper we extend the parameter space over which polynomial average time results are known. We present an algorithm called *INFREQ* and show that it has polynomial average time performance over a range of parameter values including region **II** of Figure 1; no published analysis has shown this region covered by a polynomial average time algorithm.

INFREQ uses substitution rules to eliminate or combine clauses containing in-

frequent variables: that is, variables occurring only once or twice in an instance. Infrequent variables are also eliminated by applying these rules. In addition, *INFREQ* checks the input for a null clause before processing. If one is found, *INFREQ* immediately stops with the result that the input is unsatisfiable. Otherwise, *INFREQ* does an exhaustive search over all remaining variables for a solution.

The results of this paper show that exhaustive search over the variables which are not infrequent is, for certain relationships between the parameters p , n , and r , speeded up considerably as a result of the null check and the preprocessing. The idea seems to be generalizable and may represent the first of a family of such results that will take care of a large portion of the remaining parameter space for which polynomial average time results are not now known. The result of such a generalization, to the extent that it is possible, is apparent from the analysis presented here.

Specifically, the results of this paper are that random instances of Satisfiability are “easy” in the average case if $n = r^\epsilon$, $2/3 > \epsilon > 0$, and $pr < (\ln(n)/4)^{1/3}r^{2/3-\epsilon}$; or $n = r^\epsilon$, $1 > \epsilon \geq 2/3$, and $pr \leq (1 - \epsilon - \delta)\ln(n)/\epsilon$ for any $\delta > 0$; or $pn \rightarrow 0$, and $pr < \gamma \ln \ln(n)$ for any $\gamma > 0$. These results include region II in Figure 1, a region not covered by a published polynomial average time result. The first of these results is due to the resolution component and does not depend on the presence of null clauses in an instance. Thus, in Figure 1, the left boundary of region III is due to limited resolution. As will be explained in the remarks following Theorem 1, random instances generated according to the parameter space of the first result have relatively few variables which are not infrequent. Therefore, the exhaustive search in *INFREQ* is over a sufficiently small number of variables to obtain polynomial average time. The second result depends on checking for null clauses *and* the elimination of infrequent variables. In Figure 1, the portion of region II that is below an imaginary extension of the lower boundary of region III is due to this effect. In this case the average number of variables which are not infrequent is considerable but *INFREQ* has polynomial average time because of the combination of null clause check and exhaustive search.

2. The Algorithm. For convenience, we write a clause as a tuple of the literals it contains. For example, the clause $(x \vee y \vee z)$ is written as (x, y, z) . Similarly, we write the conjunction of two clauses $C_1 \wedge C_2$ as C_1, C_2 .

Let a variable which appears exactly once in an instance I be called a *unit* variable. Let a variable which appears exactly twice in I be called a *double* variable. Let a variable which appears at least three times in I be called a *serious* variable. The table below defines substitutions for clauses in I containing unit and double variables. In the table we use v to denote a positive literal taken from a unit or double variable, \bar{v} a negative literal so taken, and x and y either a positive or negative literal which is not necessarily taken from a unit or double variable.

var type	substitution name	occurrence	replacement
unit	unit elimination	(v, x, \dots)	$true$
		(\bar{v}, x, \dots)	
double	double elimination	(v, v, x, \dots)	
		$(\bar{v}, \bar{v}, x, \dots)$	
	trivial elimination	(v, \bar{v}, x, \dots)	
	pure literal rule	$(v, x, \dots), (v, y, \dots)$	
		$(\bar{v}, x, \dots), (\bar{v}, y, \dots)$	
resolution	$(v, x, \dots), (\bar{v}, y, \dots)$	(x, \dots, y, \dots)	

When we say *apply unit elimination* we mean, according to the table above, look for

a clause containing a unit variable v and replace it with the logical value *true*; if no such clause exists do nothing. Similar statements hold for applying any of the other substitution rules listed in the table. It is possible that, after repeated applications of double-variable substitution rules, some double variables will occur only once in I . By *clean up double variables* we mean eliminate all clauses containing double variables that appear once in I .

Consider the following algorithm for solving instances of Satisfiability:

INFREQ(I) :

1. If I has a null clause then return “unsatisfiable”
2. Otherwise,
 - a. Repeatedly apply double variable substitution rules in order until opportunities vanish
 - b. Clean up all remaining double variables
 - c. Repeatedly apply unit elimination until opportunities vanish
 - d. For all truth assignments t to serious variables in I , if t satisfies I then return “satisfiable”
3. Return “unsatisfiable”

End *INFREQ*

In step (2d) *INFREQ* terminates as soon as the first satisfiable truth assignment is discovered. It should be apparent that the size of I is not increased by the application of *INFREQ* to I . It should also be apparent that all unit and double variables are eliminated from I in steps (2a), (2b), and (2c) of *INFREQ* (these are the preprocessing steps). Thus, in step (2d), the truth assignment t is an assignment to all variables which appear in the processed I .

LEMMA 2.1. *INFREQ* returns “satisfiable” if and only if I is satisfiable.

Proof. The proof is straightforward and is omitted. \square

3. The Analysis. To simplify the analysis, we show that the expected number of iterations in step (2d) of *INFREQ* is bounded by a polynomial in n under several conditions. Since the complexity of each step is polynomially bounded, the average running time of *INFREQ* must then be polynomially bounded under those conditions as well.

Let $I_=(y)$ denote the event that the input contains exactly y serious variables. Let $I_{\geq}(y)$ denote the event that the input contains at least y serious variables. Let I_{ϕ} denote the event that the input contains a null clause. Let $T(n, r, p)$ denote the average number of steps executed by *INFREQ* given that instances are generated according to model $M(n, r, p)$. Then, since the number of steps required by exhaustive search on an input with exactly y serious variables is at most 2^y , we can write

$$\begin{aligned}
T(n, r, p) &\leq Pr(I_{\phi}) + \sum_{y=1}^r 2^y \cdot Pr(\bar{I}_{\phi} \wedge I_=(y)) \\
&= Pr(I_{\phi}) + 2Pr(\bar{I}_{\phi} \wedge I_=(1)) + \sum_{y=2}^r 2^y \cdot Pr(\bar{I}_{\phi} \wedge I_=(y))
\end{aligned}$$

$$\begin{aligned}
&= Pr(I_\phi) + 2Pr(\bar{I}_\phi \wedge I_=(1)) + \sum_{y=2}^r \left(1 + \sum_{x=0}^{y-1} 2^x\right) Pr(\bar{I}_\phi \wedge I_=(y)) \\
&= Pr(I_\phi) + 2Pr(\bar{I}_\phi \wedge I_=(1)) + \sum_{y=2}^r Pr(\bar{I}_\phi \wedge I_=(y)) + \sum_{y=2}^r \sum_{x=0}^{y-1} 2^x \cdot Pr(\bar{I}_\phi \wedge I_=(y)). \\
&= Pr(I_\phi) + Pr(\bar{I}_\phi \wedge I_=(1)) + Pr(\bar{I}_\phi \wedge I_{\geq}(1)) + \sum_{y=2}^r \sum_{x=0}^{y-1} 2^x \cdot Pr(\bar{I}_\phi \wedge I_=(y)).
\end{aligned}$$

Interchanging the order of summation in the double sum gives

$$\begin{aligned}
T(n, r, p) &\leq Pr(I_\phi) + Pr(\bar{I}_\phi \wedge I_=(1)) + Pr(\bar{I}_\phi \wedge I_{\geq}(1)) + \sum_{x=1}^{r-1} \sum_{y=x+1}^r 2^x \cdot Pr(\bar{I}_\phi \wedge I_=(y)) \\
&\quad + \sum_{y=2}^r Pr(\bar{I}_\phi \wedge I_=(y)) \\
&= Pr(I_\phi) + Pr(\bar{I}_\phi \wedge I_=(1)) + \sum_{y=2}^r Pr(\bar{I}_\phi \wedge I_=(y)) + Pr(\bar{I}_\phi \wedge I_{\geq}(1)) \\
&\quad + \sum_{x=1}^{r-1} 2^x \sum_{y=x+1}^r Pr(\bar{I}_\phi \wedge I_=(y)) \\
&= Pr(I_\phi) + 2Pr(\bar{I}_\phi \wedge I_{\geq}(1)) + \sum_{x=1}^{r-1} 2^x \cdot Pr(\bar{I}_\phi \wedge I_{\geq}(x+1)) \\
&\leq 1 + 2Pr(\bar{I}_\phi \wedge I_{\geq}(1)) + \sum_{x=1}^{r-1} 2^x \cdot Pr(\bar{I}_\phi \wedge I_{\geq}(x+1)) \\
&\leq 3 + \sum_{x=1}^{\lfloor 6\mu \rfloor} 2^x \cdot Pr(\bar{I}_\phi) + \sum_{x=\lfloor 6\mu \rfloor+1}^{r-1} 2^x \cdot Pr(I_{\geq}(x+1)) \\
(1) \quad &= 3 + \sum_{x=1}^{\lfloor 6\mu \rfloor} 2^x \cdot Pr(\bar{I}_\phi) + \sum_{x=\lfloor 6\mu \rfloor+2}^r 2^{x-1} \cdot Pr(I_{\geq}(x))
\end{aligned}$$

where μ is the mean number of serious variables in an instance. The appearance of the number 6 in (1) will be explained below.

First, we obtain a bound on the second sum in (1). Since variables are placed independently in clauses, the number of serious variables in an instance is binomially distributed. By the Chernoff bound for binomial distributions [6], $Pr(I_{\geq}((1+\beta)\mu)) < e^{-\ln(1+\beta)\beta\mu/2}$, $\beta > 0$. Below we shall make use of this and the easily verified fact that $x \ln(2) - \ln(x/\mu)(x/\mu - 1)\mu/2 < 0$ if $x \geq \lfloor 6\mu \rfloor + 1$ (this is the reason why 6 appears in (1)). Thus,

$$\begin{aligned}
\sum_{x=\lfloor 6\mu \rfloor+2}^r 2^{x-1} \cdot Pr(I_{\geq}(x)) &\leq \sum_{x=\lfloor 6\mu \rfloor+2}^r 2^x e^{-\ln(x/\mu)(x/\mu-1)\mu/2} \\
&= \sum_{x=\lfloor 6\mu \rfloor+1}^r e^{x \ln(2) - \ln(x/\mu)(x-\mu)/2}
\end{aligned}$$

$$\leq \sum_{x=\lceil 6\mu \rceil + 1}^r 1 \leq r.$$

Next, we obtain an upper bound on the first sum in (1). The probability that a clause is null is $(1-p)^{2r}$. Hence, the probability that all clauses are not null is $Pr(\bar{I}_\phi) = (1 - (1-p)^{2r})^n$. Thus, we write

$$(2) \quad \sum_{x=1}^{\lfloor 6\mu \rfloor} 2^x \cdot Pr(\bar{I}_\phi) = Pr(\bar{I}_\phi) \sum_{x=1}^{\lfloor 6\mu \rfloor} 2^x = (1 - (1-p)^{2r})^n \sum_{x=1}^{\lfloor 6\mu \rfloor} 2^x.$$

It may be verified that $(1-p) \geq e^{-p-p^2}$ if $0 \leq p < 1/2$. Hence

$$(3) \quad \begin{aligned} \sum_{x=1}^{\lfloor 6\mu \rfloor} 2^x \cdot Pr(\bar{I}_\phi) &\leq (1 - (1-p)^{2r})^n 2^{\lfloor 6\mu \rfloor + 1} \\ &\leq (1 - e^{-2rp(1+p)})^n 2^{\lfloor 6\mu \rfloor + 1} \leq e^{-ne^{-2p(1+p)r}} 2^{6\mu + 1} \\ &= e^{-ne^{-2p(1+p)r} + \ln(2)(6\mu + 1)}. \end{aligned}$$

Now, we compute μ and obtain upper bounds on (3). The probability that a variable is not in a particular clause is the probability that neither literal associated with the variable is in that clause and is equal to $(1-p)^2$. Since clauses are independently chosen, the probability that a variable is not in a given instance is $(1-p)^{2n}$, the probability that a variable appears once in an instance is $2pn(1-p)^{2n-1}$, and the probability that a variable appears twice in an instance is $\binom{2n}{2}p^2(1-p)^{2n-2}$. Therefore, the probability that a variable is a serious variable is $1 - (1-p)^{2n} - 2pn(1-p)^{2n-1} - n(2n-1)p^2(1-p)^{2n-2}$ which may be reduced to $1 - (1-p)^{2n}(1 + 2pn/(1-p) + 2(pn)^2(1-1/n)/(1-p)^2)$.

THEOREM 3.1. *INFREQ runs in polynomial average time if $n = r^\epsilon$, $\epsilon \leq 2/3$, and $pr \leq (\ln(n)/4)^{1/3}r^{2/3-\epsilon}$, or if $n = r^\epsilon$, $1 > \epsilon > 2/3$, and $pr < (1 - \epsilon - \delta)\ln(n)/\epsilon$, for any $\delta > 0$, or if $pn \rightarrow 0$ and $pr < \gamma \ln \ln(n)$ for any $\gamma > 0$.*

Proof. If $n = r^\epsilon$, $\epsilon < 2/3$, and $pr \leq (\ln(n)/4)^{1/3}r^{2/3-\epsilon}$ then $pn = pr \cdot r^{\epsilon-1} \leq (\ln(n)/4)^{1/3}r^{2/3-\epsilon+\epsilon-1} = (\ln(n)/4)^{1/3}r^{-1/3} \rightarrow 0$. If $1 > \epsilon \geq 2/3$ and $pr \leq \ln(n)$ then $pn = pr \cdot r^{\epsilon-1} \leq \epsilon \ln(r)r^{\epsilon-1} \rightarrow 0$. So, we assume $pn \rightarrow 0$. This implies $p < 1/2$. Then the probability that a variable is serious is $1 - (1-p)^{2n}(1 + 2pn/(1-p) + 2(pn)^2(1-1/n)/(1-p)^2) = (4/3)(np)^3 + O((np)^4)$. From now on we ignore the small term for simplicity. Since variables are placed independently in clauses, the number of serious variables in an instance is binomially distributed with parameters r and $(4/3)(np)^3$. Thus, the mean number of serious variables in an instance is $\mu = (4/3)(np)^3r$. Substituting into (3) gives $e^{-ne^{-2p(1+p)r} + \ln(2)(8(np)^3r + 1)}$ which is polynomially bounded if $e^{-ne^{-2p(1+p)r} + \ln(2)(8(np)^3r)}$ is. Therefore, we require

$$(4) \quad -e^{-2p(1+p)r} + 5.545(pn)^2(pr) \leq \ln(n)/n.$$

Let $n = r^\epsilon$, $1 > \epsilon > 0$. Then, after rearranging, (4) becomes

$$(5) \quad \begin{aligned} 5.545(pr)^3 r^{2(\epsilon-1)} &\leq e^{-2pr(1+p)} + \epsilon r^{-\epsilon} \ln(r) \iff \\ 5.545(pr)^3 r^{3\epsilon-2} &\leq r^\epsilon e^{-2pr(1+p)} + \epsilon \ln(r) \end{aligned}$$

Let $\epsilon \leq 2/3$ and $pr = (\ln(n)/4)^{1/3}r^\alpha$, α a constant. Then (5) becomes

$$(6) \quad 5.545 \ln(n) r^{3\alpha+3\epsilon-2}/4 \leq r^\epsilon e^{-2r^\alpha-2r^{2\alpha-1}} + \epsilon \ln(r).$$

Clearly, (6) is satisfied if $3\alpha + 3\epsilon - 2 \leq 0$ or $\alpha \leq 2/3 - \epsilon$. Thus, if $\epsilon \leq 2/3$ and $pr < (\ln(n)/4)^{1/3}r^{2/3-\epsilon}$ then (2) is polynomially bounded.

Now let $1 > \epsilon > 2/3$ and $pr = \alpha \ln(n) = \alpha \epsilon \ln(r)$. Then (5) becomes

$$(7) \quad 5.545(\alpha \ln(r))^3 r^{3\epsilon-2} \leq r^{\epsilon-2\alpha\epsilon(1+\ln(n)/r)} + \epsilon \ln(r).$$

Inequality (7) is satisfied if $3\epsilon - 2 \leq \epsilon - 2\alpha\epsilon - 2\delta$ for any positive constant δ and this is satisfied if $\alpha \leq (1 - \epsilon - \delta)/\epsilon$. Thus, if $1 > \epsilon > 2/3$ and $pr \leq (1 - \epsilon - \delta) \ln(n)/\epsilon$ then (2) is polynomially bounded.

The remaining case, $pn \rightarrow 0$ and $pr < \gamma \ln \ln(n)$, is straightforward. \square

We make four remarks about the proof of Theorem 1. First, in (6) only the term $e^{-2r^\alpha-2r^{2\alpha-1}}r^\epsilon$ is due to the presence of null clauses in I . But this term is ignored when determining that $\alpha \leq 2/3 - \epsilon$ in the sentence following (6). Thus, the polynomial average time result for $n = r^\epsilon$, $\epsilon < 2/3$, is *not* due to the presence of null clauses in I .

Second, in (7) the term $r^{\epsilon-2\alpha\epsilon(1+\ln(n)/r)}$ is due to the presence of null clauses and the term $r^{3\epsilon-2}$ is due to the removal of infrequent variables. Both the null clause check and removal of infrequent variables account for polynomial average time when $pr < (1 - \epsilon) \ln(n)/\epsilon$, $2/3 < \epsilon < 1$. That is, neither checking for null clauses alone nor removing infrequent variables without checking for null clauses is powerful enough to achieve this result.

Third, if $pr < (\ln(n)/4)^{1/3}r^{2/3-\epsilon}$, $0 < \epsilon < 2/3$, then the average number of serious variables, although possibly an increasing function of n , is small compared to the number of infrequent variables. Thus, in this case it can be said that *INFREQ* works well because nearly all the variables are eliminated by resolution, the pure literal rule, unit elimination, double elimination, and trivial elimination. On the other hand, if $pr < (1 - \epsilon) \ln(n)/\epsilon$ and $2/3 < \epsilon < 1$ then it can turn out that many variables are serious. In fact, if $pr = (1 - \epsilon) \ln(n)/\epsilon$, $2/3 < \epsilon < 1$, then the average number of serious variables is $\theta(\ln^3(r)r^{3\epsilon-2})$. Exhaustive search over so many variables would require superpolynomial time. However, *INFREQ* works well on the average in this case too.

Removing resolution (recall this is on double variables only) and the null clause check from *INFREQ* leaves essentially the algorithm that was analysed in [1] called *PLR*. But *PLR* requires superpolynomial average time if $pr < (1 - \epsilon) \ln(n)/\epsilon$, $2/3 < \epsilon < 1$. Thus, resolution and the null check account for the good average performance of *INFREQ* in this case. Moreover, *PLR* requires superpolynomial average time even if $pr > \sqrt{\ln(n)} \cdot r^{1/2-\epsilon}$, $1/2 < \epsilon < 2/3$. Thus, in the case where $\sqrt{\ln(n)} \cdot r^{1/2-\epsilon} < pr < (\ln(n)/4)^{1/3}r^{2/3-\epsilon}$, $1/2 < \epsilon < 2/3$, only the addition of resolution on double variables to *PLR* accounts for polynomial average time. This means that with $1/2 < \epsilon < 2/3$ large samples of instances with up to r^α literals per clause on the average, $0 < \alpha < 1/6$, can be solved in polynomial average time with *INFREQ* whereas with *PLR* superpolynomial average time is required even if the average number of literals per clause is vanishingly small. It is perhaps surprising that such a small change to *PLR* can have such an effect on average case performance.

The fourth remark concerns the scope of infrequent variables. From the paragraph preceding Theorem 1 it should be evident that, if $pn \rightarrow 0$ and only unit variables are

eliminated in *INFREQ*, then $\mu = \theta((np)^2 r)$ and, up to constant factors, (6) becomes

$$\ln(n)r^{2\alpha+2\epsilon-1} \leq r^\epsilon e^{-2r^\alpha - 2r^{2\alpha-1}} + \epsilon \ln(r)$$

which is satisfied if $\alpha < 1/2 - \epsilon$. If we could use substitution rules to eliminate triple variables, those which appear three times in an instance, then $\mu = \theta((np)^4 r)$ and $\alpha < (3/4) - \epsilon$. If we could eliminate all variables occurring i or fewer times in I then we would have $\mu = \theta((np)^{(i-1)} r)$ and polynomial average time if $pr < r^{i/(i+1)-\epsilon}$, $\epsilon < i/(i+1)$. Clearly i does not have to be very large to make a major impact on the parameter space supporting polynomial average time. Unfortunately, trying to eliminate even triple variables can cause an exponential explosion of the size of I . In this event the assumption that the complexity of each step of *INFREQ* is polynomially bounded is not valid. We ask: are explosions so infrequent that they do not significantly affect average time performance? An affirmative answer would have a major impact on polynomial average time results under the random-clause-size model. We leave investigation of this question for a future paper.

The next theorem shows where *INFREQ* runs in polynomial average time when $n = \beta r$, β a positive constant.

THEOREM 3.2. *INFREQ runs in polynomial average time if $n/r = \beta$, where β is a positive constant, and $2.64(1 - (1-p)^{2\beta r}(1 + 2\beta pr + 2(\beta pr)^2)) < \beta e^{-2pr}$.*

Proof. Since $p < 1$, $1/(1-p) > 1$ and $1/(1-p)^2 > 1$. Then

$$\mu = (1 - (1-p)^{2n}(1 + 2pn/(1-p) + 2(pn)^2/(1-p)^2))r \leq (1 - (1-p)^{2n}(1 + 2pn + 2(pn)^2))r.$$

Thus, (3) is polynomially bounded if

$$(8) \quad \begin{aligned} -ne^{-2pr} + \ln(2)(6(1 - (1-p)^{2n}(1 + 2pn + 2(pn)^2)))r &\leq \ln(n) \iff \\ -\beta e^{-2pr} + 4.15(1 - (1-p)^{2\beta r}(1 + 2\beta pr + 2(\beta pr)^2)) &\leq \ln(n)/r. \end{aligned}$$

The theorem follows. \square

According to Theorem 2, *INFREQ* has polynomial average time if $2pr < \ln(\beta) - \ln(2.64)$ (this is fairly tight if β is large). If $\beta = 1$ then *INFREQ* has polynomial average time if $pr < .5$.

4. Conclusions. We have investigated a simple strategy for solving instances of CNF Satisfiability with respect to average case performance. The important idea is the elimination of infrequent variables before applying, in this case, exhaustive search. We have shown that this strategy is superior in average case performance to all other algorithms analyzed under the random-clause-size model when $pr < \sqrt{\epsilon r \ln(r)}$, $n < r^\epsilon$, and $\epsilon < 1$. The strategy may be generalizable, to some extent, and the analysis seems to suggest the outcome of an investigation of such a generalization.

REFERENCES

- [1] Bugrara, K., Pan, Y., and Purdom, P. W., "Exponential average time for the pure literal rule", *SIAM J. Comput.*, Vol. 18, No. 2, (1989) pp. 409-418.
- [2] Franco, J., "On the probabilistic performance of algorithms for the Satisfiability problem", *Information Processing Letters*, Vol. 23, (1986) pp. 103-106.
- [3] Iwama, K., "CNF Satisfiability test by counting and polynomial average time", *SIAM J. Comput.*, Vol. 18, No. 2, (1989) pp. 385-391.

- [4] **Purdom, P. W., and Brown, C. A.**, "The pure literal rule and polynomial average time", *SIAM J. Comput.*, Vol. 14, No. 4, (1985) pp. 943-953.
- [5] **Purdom, P. W., and Brown, C. A.**, "Polynomial average time Satisfiability problems", *Information Sciences*, Vol. 41, (1987) pp. 23-42.
- [6] **Purdom, P. W., and Brown, C. A.**, "The Analysis of Algorithms", Holt, Rinehart, Winston (1985).