# The 3-XORSAT threshold

## Olivier Dubois, Jacques Mandler

**LIP6, Box 169, CNRS-Université Paris 6, 4 place Jussieu, 75252 Paris cedex 05, France**

**Note presented by Michel Talagrand.**

**Abstract**
We prove the existence of a threshold phenomenon regarding the random 3-XORSAT problem (or more generally $k$-XORSAT). We provide the value of the threshold as the solution of two transcendental equations. These results confirm rigorously those obtained by physicists using the one-step replica symmetry breaking method and thus give for the first time the proof of the validity of this method for a problem of the class of satisfiability problems. *To cite this article: O. Dubois, J. Mandler, C. R. Acad. Sci. Paris, Ser. I 335 (2002) 963–966.*
© 2002 Académie des sciences/Éditions scientifiques et médicales Elsevier SAS


### Le seuil de 3-XORSAT

**Résumé**
Nous démontrons l'existence d'un phénomène de seuil pour le problème 3-XORSAT aléatoire (et plus généralement $k$-XORSAT). Nous fournissons la valeur du seuil comme solution de deux équations transcendantes. Ces résultats confirment rigoureusement ceux obtenus par des physiciens au moyen de la méthode des répliques à un pas de brisure de symétrie et apportent ainsi pour la première fois une preuve de la validité de la méthode des répliques avec brisure sur un problème de la classe des problèmes de satisfaisabilité. *Pour citer cet article : O. Dubois, J. Mandler, C. R. Acad. Sci. Paris, Ser. I 335 (2002) 963–966.*
© 2002 Académie des sciences/Éditions scientifiques et médicales Elsevier SAS

**1.** The 3-XORSAT (more generally, $k$-XORSAT) problem is a variant of the satisfiability problem in which a clause of length 3 is said to be satisfied iff the exclusive-OR of the values of its literals is 1. Using the fact that the exclusive-OR is just addition modulo 2 and that complementing is just adding 1 modulo 2, it is easy to see that a 3-XORSAT formula can be identified with a *system of linear equations on* $\mathbb{F}_2$, each equation having exactly 3 variables.

Let $\Omega_{m,n}$ be the set of 3-XORSAT formulae over a set of $n$ variables (not necessarily all present within each formula), comprising exactly $m$ clauses with $m = cn(1 + o(1))$ as $n \to \infty$. Experiments suggest the existence of a *satisfiability threshold* $c_0$, close to 0.92, such that for $c < c_0$ [resp. $c > c_0$], almost all formulae in $\Omega_{m,n}$ are satisfiable [resp. unsatisfiable] [2]. Let $\Psi_{m,n}$ be the set of 3-XORSAT formulae with $m$ equations in $n$ variables, having each *at least two* occurrences, endowed with the uniform probability distribution. We first show that satisfiability in $\Psi_{m,n}$ has, with respect to the parameter $c \sim m/n$, a threshold located exactly at $c = 1$. In the precise model which we adopt, a formula in $\Psi_{m,n}$ is considered to be ordered and each clause in a formula is also considered to be ordered.

THEOREM 1. – *For $m = (1 + o(1))cn$ with $c < 1$ [resp. $c > 1$], almost all formulae in $\Psi_{m,n}$ are satisfiable [resp. unsatisfiable].*

To show that for $c > 1$ almost all formulae in $\Psi_{m,n}$ are unsatisfiable is a straightforward application of the first moment method, since the expected number $\mathbf{E}(N)$ of solutions of a random formula $\omega \in \Psi_{m,n}$ is $2^{m-n}$ (as in the case of $\Omega_{m,n}$). To show that for $c < 1$ almost all formulae in $\Psi_{m,n}$ are satisfiable is a much more involved application of the second moment inequality: $\Pr(N > 0) \geqslant \mathbf{E}(N)^2/\mathbf{E}(N^2)$. The second moment is the quotient of the number of formulae satisfied by a given pair of truth assignments to the variables, summed over the assignment pairs, to the total number of formulae. We take as parameters the proportion $\alpha$ of variables having the same value in both assignments, and the proportion $r$ of the $3m$ places in a random ordered formula which receive one of these $\alpha n$ variables. The r.h.s. of the system is irrelevant, providing we impose as a compatibilty condition that each equation comprises either 3 variables with the same value in both assignments, i.e., 3 of the $\alpha n$ variables, or 2 variables each with different values in both assignments, i.e., 2 of the $(1 - \alpha)n$ variables, and a third one with the same value, i.e., 1 of the $\alpha n$ variables. Setting $I_n = \{0, 1/n, 2/n, \ldots, 1\}$ this leads to

$$\mathbf{E}(\mathbf{N}^2) = \#\mathcal{F}^{-1} \sum_{\alpha \in I_n} \sum_{r \in [1/3,1] \cap I_{3m}} 2^n \binom{n}{\alpha n} \binom{m}{(1-r)3m/2} 3^{(1-r)3m/2} S(r3m, \alpha n, 2) \alpha n!$$
$$\times S\big((1-r)3m, (1-\alpha)n, 2\big) \big[(1-\alpha)n\big]!$$

where $S(p, q, 2)$ is the level-2 generalized Stirling number of the second kind, or the number of partitions of $q$ objects into $p$ subsets having each at least two elements, and $\#\mathcal{F}$ is the total number of formulae in $\Psi_{m,n}$, that is : $\#\mathcal{F} = 2^m S(3m, n, 2)n!$. Using an asymptotic estimate for $S(p, q, 2)$ which has been established in [6], we can show that the exponential order of $\mathbf{E}(\mathbf{N}^2)$ is bounded by $\{\max_{\alpha \in [0,1], r \in [1/3,1]} \exp[f(\alpha, r)]\}^n$, with

$$f(\alpha, r) = (1-c)\ln 2 - \alpha \ln \alpha - (1-\alpha)\ln(1-\alpha) + \alpha \ln(e^{x_2} - 1 - x_2) - 3rc \ln x_2$$
$$+ (1-\alpha)\ln(e^{x_1} - 1 - x_1) - 3c(1-r)\ln x_1 + 3rc \ln r + 3c(1-r)\ln(1-r) + 3c \ln x_0$$
$$- 3/2c(1-r)\ln[1/2(1-r)] - c(1 - 3/2(1-r))\ln(1 - 3/2(1-r)) - \ln(e^{x_0} - 1 - x_0)$$

the $x_i$ being given implicitly by $(q_i/p_i)x_i = (e^{x_i} - 1 - x_i)/(e^{x_i} - 1)$, where $p_0 = 3m$, $q_0 = n$, $p_1 = r3m$, $q_1 = \alpha n$, $p_2 = (1 - r)3m$, and $q_2 = (1 - \alpha)n$. Analytical and algebraic manipulations allow an explicit investigation of all stationary points of $f$, with the conclusion that the only local maximum is at $(1/2, 1/2)$ where $f$ evaluates to $2(1 - c)\ln 2$. Knowing this, the ratio $\mathbf{E}(N^2)/\mathbf{E}(N)^2$, written in the form $\mathbf{E}(N^2)/\mathbf{E}(N)^2 \sim (1/n)\sum_{j=0}^n \sum_{l=m}^{3m} g(j/n, l/(3cn))\exp(n\,h(j/n, l/(3cn)))$, where $h(\alpha, r) = f(\alpha, r) - 2(1 - c)\ln 2$, can be studied by a discrete bidimensional Laplace method, giving $\frac{\mathbf{E}(N^2)}{\mathbf{E}(N)^2} \sim g(1/2, 1/2)\frac{3c\pi}{\sqrt{D}}$. Here $D$ is the determinant of the Hessian matrix of $h$ at $(1/2, 1/2)$, found to equal $16x_0/(1 - R)$ with $R = x_0 e^{x_0}/(e^{x_0} - 1)^2$, while using exact asymptotic equivalents for all quantities involved yields $g(1/2, 1/2) = 4/\pi[x_0(3c - 1) - 3c(3c - 2)]^{-1/2}$. It follows from the expressions of $D$ and $g(1/2, 1/2)$, and the equation for $x_0$, that $D = (3c\pi)^2 g(1/2, 1/2)^2$, so that $\mathbf{E}(N^2)/\mathbf{E}(N)^2$ equals 1 in the limit, giving the required lower bound of 1 for the probability of satisfiability.

**2.** In a combinatorial approach, we can then analyze the following algorithm (with a satisfiability-preserving main loop) acting on a random formula $\omega \in \Omega_{m,n}$.

ALGORITHM A:
**while** $\omega$ contains literals with a single occurrence **do**
Let $\pi = \{$All single-occurrence literals in $\omega\}$
Remove from $\omega$ all clauses containing a literal from $\pi$;
**od**
**if** $\omega = \emptyset$ **or** the formula has fewer clauses than effectively present variables **then Success else Failure**.

To study Algorithm A, let us denote by $\Phi_{m,n,p}$ the (uniformly probabilized) space of formulae with $m$ equations in $n$ effectively present variables, $p$ of which have a single occurrence. Further, in a random formula $F$ from $\Omega_{n,m}$, let the r.v.'s $\bar{n}$ and $\bar{p}$ denote, respectively, the numbers of effectively present, and single-occurrence variables. We show that, conditionally on the values of $\bar{n}$ and $\bar{p}$, $F$ is uniformly distributed in $\Phi_{m,\bar{n},\bar{p}}$, and then, for any fixed $\delta \in ]1/2, 1[$, by a large-deviation analysis, that q.s. (quite surely, see [1]): $\bar{n} = n(1 - \exp(-3m/n)) + \mathrm{O}(n^\delta)$, $\bar{p} = 3m \exp(-3m/n) + \mathrm{O}(n^\delta)$. Also, we prove *maintenance of uniformity* at each step: if an iteration of Algorithm A sends a random formula of $\Phi_{m,n,p}$ to a formula in $\Phi_{m',n',p'}$, then the latter is also random, conditionally on the values of $m'$, $n'$ and $p'$.

Regard, then, the $i$th iteration of Algorithm A as sending a random element of $\Phi_{m_i,n_i,p_i}$ to one of $\Phi_{m_{i+1},n_{i+1},p_{i+1}}$, with $m_1 = m$, $n_1 = \bar{n}$, $p_1 = \bar{p}$. Let $c_i$ be the q.s. limit of $m_i/n_i$. Our main result is then:

THEOREM 2. – *Define $f_c(x) = 1/(3(1 - \sqrt{3c/x}(e^{-x} - 1 + x)/x))$ and $g_c(\lambda) = 3c(1 - e^{-\lambda})^2$. Then, for any $i \geqslant 1$, $c_i = f_c(g_c^{(i-1)}(3c))$, where the superscript denotes iterated composition.*

Theorem 2, together with Theorem 1, implies the existence of the 3-XORSAT threshold and gives its value. Indeed, since $g_c(3c) < 3c$, the iterates $g_c^{(i-1)}(3c)$ tend to the largest fixed-point $\tilde{\lambda}_c$ of $g_c$. There is exactly one value $c_0$ of $c$ such that $f_c(\tilde{\lambda}_c) = 1$. (It so happens that $c_0$ is also the unique $c$ such that $\max_x f_c(x) = 1$.) For $c < c_0$, the maximum of $f$ is $< 1$ anyway, so for large $n$ Algorithm A almost surely answers **Success**, and a.e. formula in $\Omega_{n,m}$ is satisfiable. For $c > c_0$, $f_c(\tilde{\lambda}_c) > 1$, so similarly a.e. formula is unsatisfiable. Thus:

THEOREM 3. – *The 3-XORSAT satisfiability threshold is the value of $c$ in the unique solution $(c_0, \tilde{\lambda})$ with $\tilde{\lambda} \neq 0$ to the system $g_c(\lambda) = \lambda$, $f_c(\lambda) = 1$.*

Up to the change of unknown $u = \sqrt{\lambda/(3c)}$, this system is the same as that in [3], giving absolute confirmation of the validity of the replica symmetry breaking method as applied to the XORSAT problem.

To prove Theorem 2, we first analyze the general step as in [1]. Using a generalized form of Poissonization, occurrences of variables in formulae of $\Phi_{m,n,p}$ are seen to be 'concentrated around left-truncated Poisson means' in the sense that the number $\nu_k$ of variables having $k$ occurrences in a random formula $\omega \in \Phi_{m,n,p}$ is, for $2 \leqslant k \leqslant \log n$, almost surely given by $\nu_k = (n - p)\hat{\lambda}^k/((e^{\hat{\lambda}} - 1 - \hat{\lambda})k!) + \mathrm{O}(n^\delta)$, with $\hat{\lambda}$ defined by $\hat{\lambda}(e^{\hat{\lambda}} - 1)/(e^{\hat{\lambda}} - 1 - \hat{\lambda}) = (3m - p)/(n - p)$. We combine this with concentration inequalities and counting arguments, to the effect that

THEOREM 4. – *Suppose that $\omega$ is chosen uniformly from $\Phi_{m,n,p}$, with $m, p > n^\delta$, and all clauses containing single-occurrence variables are deleted, giving the formula $\omega' \in \Phi_{m',n',p'}$. Then q.s. $m' = m(1 - \alpha)^3 + \mathrm{O}(n^\delta)$, $n' = (n - p)(1 - \beta) + \mathrm{O}(n^\delta)$, $p' = (n - p)\gamma + \mathrm{O}(n^\delta)$, where $\alpha = p/(3m)$ is the probability that a random literal in $\omega$ occurs only once, and where, denoting $\tau = \exp((2\alpha - \alpha^2)\hat{\lambda}) - 1$, we set $\beta = 1/(e^{\hat{\lambda}} - 1 - \hat{\lambda})(\tau - (2\alpha - \alpha^2)\hat{\lambda})$ and $\gamma = \hat{\lambda}/(e^{\hat{\lambda}} - 1 - \hat{\lambda})(1 - \alpha)^2\tau$.*

Denoting by $\hat{\lambda}_i$ the value of $\hat{\lambda}$ (in the large $n$ limit) at step $i$, we then show, by induction on $i$, that the large-$n$ q.s. limits of $m_i/n$, $n_i/n$, and $p_i/n$ are $\hat{\lambda}_i(1 - e^{-\hat{\lambda}_{i-1}})/3$, $1 - e^{-\hat{\lambda}_i} - \hat{\lambda}_i e^{\hat{\lambda}_{i-1}}$, and $\hat{\lambda}_i(e^{-\hat{\lambda}_i} - e^{-\hat{\lambda}_{i-1}})$, respectively, and in doing so we obtain the relation $\hat{\lambda}_{i+1} = (1 - \alpha_i)^2\hat{\lambda}_i$ with $\alpha_i$ the limit of $p_i/(3m_i)$. Observing that $1 - \alpha_i = (1 - e^{-\hat{\lambda}_i})/(1 - e^{-\hat{\lambda}_{i-1}})$ and therefore that $\hat{\lambda}_{i+1}/(1 - e^{-\hat{\lambda}_i})^2$ is a constant, Theorem 2 follows.

**3**. Additionally, as suggested by Monasson [4], we show that the analytical equations of Theorem 3 can also be derived by differential calculus methods, using a theorem of Wormald [7]. We obtain the equations in $c, u$ version, namely $1 - u = \exp(-3cu^2)$, $cu^3 = u - 3cu^2(1 - u)$. We analyze a modified version of Algorithm A, where at each step a single clause is picked randomly among those containing a single-occurrence variable, then removed. We follow the evolution in 'time', $T$, of the random variables $n_j(T)$ and $m(T)$ equal, respectively, to the number of variables with $j$ occurrences and the number of clauses at

iteration $T$. We focus on $n_0$ and $n_1$. The total number of variables, $n$, is now considered a constant, while $n_0$ continually increases. Using the concentration of the $n_j$'s ($j \geqslant 2$) around truncated Poisson means, it is observed that, conditional on the current state, the expected variations of $n_0$, $n_1$, and $m$ during the $(T+1)$st step of the algorithm are $f_i(T/n, n_0/n, n_1/n, m/n)$ for $i = 0, 1, 2$, repectively, where, on a suitable domain $D \subset \mathbb{R}^4$,

$$f_0(t, y_0, y_1, y_2) = 1 + \frac{2y_1}{3y_2},$$

$$f_1(t, y_0, y_1, y_2) = -1 - \frac{2y_1}{3y_2} + 2\left(1 - \frac{y_1}{3y_2}\right)\frac{\phi(t, y_0, y_1, y_2)}{\exp(\phi(t, y_0, y_1, y_2)) - 1},$$

$$f_2(t, y_0, y_1, y_2) = -1,$$

$\phi$ being defined implicitly on $D$ by: $(\phi(t, y_0, y_1, y_2))^{-1} - (e^{\phi(t, y_0, y_1, y_2)} - 1)^{-1} = (1 - y_0 - y_1)/(3y_2 - y_1)$. The $f_i$'s are seen to be Lipschitz on $D$, so by Wormald's theorem, almost surely and uniformly for $0 \leqslant T < T_s$ ($T_s$ being the stopping time of the algorithm, when $n_1$ becomes 0), $n_0(T)$, $n_1(T)$ and $m(T)$ are equal to $(z_j(T/n) + o(1))n$, with respectively $j = 0, 1, 2$; and $z_j(t)$ being the unique solution in $D$ of the differential system $\{dz_j/dt = f_j(t, z_0, z_1, z_2)\}$ subject to the initial conditions $z_j(0) = n_j(0)/n + o(1) = e^{-3c}(3c)^j/j!$, for $j = 0, 1$, and $z_2(0) = m(0)/n = c$. Setting $\varpi(t) = c - t$, $u(t) = [\varpi(t)/c]^{1/3} = (1 - t/c)^{1/3}$, and $\gamma(t) = 3cu(t)^2$, we show, by direct substitution, that the solution of this differential system is $z_2 = \varpi$, $z_1 = \gamma(u - 1 + e^{-\gamma})$, $z_0 = (1 + \gamma)e^{-\gamma} - z_1$, and that the largest semi-open interval to which this solution can be extended is $[0, t_s[$, with $t_s$ the uniform limit of $T_s/n$ as $n \to \infty$. Since $n_1(T_s - 1) = 1$, we then have $\lim_{n \to \infty} z_1(T_s/n - 1/n) = \lim_{n \to \infty} 1/n = 0$, which for the values of $c$ of interest implies $1 - u(t_s) - e^{-3cu(t_s)^2} = 0$. Additionally, from the threshold result for $\Psi_{m(T_s), n - n_0(T_s)}$ and the satisfiability-preserving loop in the algorithm, it is seen that the unique real number $c_0$ below [resp. above] which the final $c = m(T_s)/[n - n_0(T_s)]$ is $< 1$ [resp. $> 1$], is indeed a threshold for 3-XORSAT, and that further $\varpi(t_s) - 1 + (1 + \gamma(t_s))e^{-\gamma(t_s)}$ is $\leqslant 0$ or $\geqslant 0$, according as $c < c_0$ or $c > c_0$. This means that, setting $t_0 = \inf_{c<c_0} t_s = \sup_{c>c_0} t_s$, the unique non-trivial solution of the above equations in $c, u$ consists of $c_0$ and $u_0 = u(t_0)$. (The expressions obtained here agree exactly with those in **2**.)

In conclusion, work like this one underlines the importance of foundational studies into the replica method, such as those of M. Talagrand (see his forthcoming book [5]).

## References

[1] A.Z. Broder, A.M. Frieze, E. Upfal, On the satisfiability and maximum satisfiability of random 3-CNF formulas, in: Proc. 4th ACM-SIAM Symp. on Discrete Algorithms, Austin, TX, 1993, pp. 322–330.

[2] N. Creignou, H. Daudé, O. Dubois, Approximating the satisfiability threshold for random $k$-XOR-formulas, 2001, submitted.

[3] S. Franz, M. Leone, F. Ricci-Tersenghi, R. Zecchina, Phys. Rev. Lett. 87 (2001) 12713–127209.

[4] R. Monasson, personal communication.

[5] M. Talagrand, Spin Glasses: A Challenge for Mathematicians, in preparation.

[6] N.M. Temme, Asymptotic estimates of Stirling numbers, Studies Appl. Math. 89 (1993) 233–243.

[7] N.C. Wormald, Differential equations for random processes and random graphs, Ann. Appl. Probab. 5 (1995) 1217–1235.