

Some Observations for the Pigeon Hole Principle

F.J. Brandenburg, J. van Leeuwen and R.B. Tan

RUU-CS-88-39
December 1988



Rijksuniversiteit Utrecht

Vakgroep informatica

Padualaan 14 3584 CH Utrecht
Corr. adres: Postbus 80.089, 3508 TB Utrecht
Telefoon 030-531454
The Netherlands

Some Observations for the Pigeon Hole Principle

F.J. Brandenburg, J. van Leeuwen and R.B. Tan

Technical Report RUU-CS-88-39
December 1988

**Department of Computer Science
University of Utrecht
P.O.Box 80.089, 3508 TB Utrecht
The Netherlands**

SOME OBSERVATIONS FOR THE PIGEON HOLE PRINCIPLE

F.J. Brandenburg* J. van Leeuwen† R.B. Tan‡

Abstract. Several elementary combinatorial facts and complexity results are proved, related to the pigeon hole principle.

Keywords and phrases: combinatorics, pigeon hole principle, subset sums, NP-completeness.

1 Introduction.

The “pigeon hole principle” (see e.g. Liu [3]) is a simple proof technique with often amazing consequences when applied to combinatorial problems. In its most elementary form, the pigeon hole principle asserts the following (for any $n \geq 1$):

“when $n + 1$ items are divided over n boxes, then at least one box receives at least 2 items”.

There are many intriguing applications of the principle in mathematical contexts, with Dirichlet’s result on the approximation of real numbers by rational numbers as perhaps the earliest example known (cf. [2]). In this note we present some additional combinatorial facts “when $n + 1$ (or more) items are divided over n boxes”, which may serve as interesting classroom illustrations or exercises when the pigeon hole principle is discussed. In section 2 we consider some aspects of the distribution of N items over n boxes, and in section 3 we prove some (new) complexity results for distributions over two sets of distinct boxes. Section 3 assumes some familiarity with the theory of NP-completeness (cf. [1]).

*Fakultät f. Mathematik u. Informatik, Universität Passau, Innstrasse 33, Postfach 2540, D-8390 Passau, Federal Republic of Germany.

†Department of Computer Science, University of Utrecht, P.O.Box 80.089, 3508 TB Utrecht, the Netherlands.

‡Department of Mathematics and Computer Science, University of Sciences and Arts of Oklahoma, Chickasha, Oklahoma 73018, USA.

2 Items in boxes.

When $n + 1$ items are divided over n boxes, it is very likely that some (many) boxes remain empty. Empty boxes play a crucial role in the proof of the following fact (for any $n \geq 1$).

Theorem A. Suppose N items are divided over n boxes, with $1 \leq N < 2n$. Then either all items end up in one box, or there is a set of k boxes containing precisely k items (for some $k \geq 1$).

Proof.

(By omitting some empty boxes it is straightforward to reduce the proof to the case in which we have $n < N < 2n$, but this is not essential for the argument.) As $N < 2n$, the boxes cannot all contain at least 2 items. If there is a box with precisely 1 item, then the theorem is proved (take $k = 1$). If there are no boxes with precisely 1 item, then we argue as follows. Let $e \geq 1$ be the number of empty boxes, let $b \geq 2$ be the smallest number of elements in any non-empty box, and assume that not all items end up in one box (which means that $b \leq n$). It follows that $(n - e)b \leq N$, hence $e \geq n - \frac{N}{b} > n - \frac{2n}{b} = (b - 2) \cdot \frac{n}{b} \geq b - 2$ and thus $e \geq b - 1$. The theorem now follows with $k = b$, by taking a box with b items and $b - 1$ boxes with zero items.

□

In the proof of theorem A, empty boxes contribute to the count in exactly the right way (i.e., they help in bringing the number of boxes up to k without changing the number of items). In the following result this effect is much less apparent, thus making it more intriguing and less straightforward than theorem A. The proof actually uses the pigeon hole principle itself!

Theorem B. Suppose N items are divided over n boxes, with $N \neq n$ and $1 \leq N \leq n + 1$. Then either all items end up in one box, or there is a set of k boxes containing precisely $N - k$ items (for some $k \geq 1$).

Proof.

Suppose N items are divided over n boxes, and suppose also that the items do not all end up in one box. It means that $n > 1$, and that there are at least two non-empty boxes. Assume without loss of generality that the 1st and the n^{th} box are non-empty. If one of these boxes contains $N - 1$ items, then the theorem follows with $k = 1$. Thus assume that all boxes contain less than $N - 1$ items.

First consider the special case that $N < n$, with each item in a distinct box (i.e., the N items end up in exactly N boxes with 1 item per box). If N is even, then the theorem follows with $k = \frac{1}{2}N$ by taking $\frac{1}{2}N$ boxes with 1 item each. If N is odd, then the theorem follows with $k = \lceil \frac{1}{2}N \rceil = \lfloor \frac{1}{2}N \rfloor + 1$ by taking $\lfloor \frac{1}{2}N \rfloor$ boxes with 1 item each and 1 box with zero items (which exists because $N < n$). Thus we may assume $N = n + 1$, or $N < n$ and the N items actually end up in at most

$N - 1$ boxes. In the latter case we simply eliminate some empty boxes from further consideration and consider only some set of $N - 1$ boxes which together contain the N items. It follows that we can assume that $N = n + 1$, for the remainder of the argument.

Add one additional item to every box. (For simplicity we refer to the additional items as “green items”, when necessary.) Let a_i denote the number of items in the i^{th} box. Then $a_i, a_n \geq 2, \sum_1^n a_i = 2n + 1$, and $1 \leq a_i \leq n$ for every i . By an application of the pigeon hole principle (cf. Liu [3], exerc. 4.47 or theorem E below) it follows that there are an $i \geq 1$ and $l \geq 0$ such that $a_i + \cdots + a_{i+l} = n$ or $a_i + \cdots + a_{i+l} = 2n$. Consider the latter case first. It is impossible that $i = 1$ and $l = n - 1$. Thus either a_1 or a_n is excluded from the sum and necessarily $a_i + \cdots + a_{i+l} \leq (2n + 1) - 2 = 2n - 1$, a contradiction. Thus we necessarily have $a_i + \cdots + a_{i+l} = n$ which, after pulling out the green items, implies that the boxes i through $i + l$ contain precisely $n - (l + 1)$ items. Thus the remaining $k = n - (l + 1)$ boxes contain precisely $(n + 1) - k = N - k$ items. (Note that $l < n - 1$, hence $k \geq 1$.) \square

From the proof it follows that theorem B holds for $N = n$ if and only if n is even. Theorem B fails for $N = n + 2$ as well: consider the distribution with 3 items in the first box, and 1 item in every other box. Then k boxes contain precisely k or $k + 2$ items, which is of opposite parity from $N - k$ and thus not equal to it for all odd N (i.e., odd n).

In an equivalent formulation for the case $N = n + 1$, theorem B asserts that either all items end up in one box or there is a set of $j = (n - k)$ boxes containing precisely $(n - j) = k$ items, for some $j \geq 1$. As either $(n + 1) - k > k$ or $k > (n - k)$, theorem B shows that the pigeon hole principle can be applied recursively within the given set of pigeon holes as well. Here is another interesting consequence.

Corollary C. Suppose $N + n$ items are divided over n boxes, with $N \neq n$ and $1 \leq N \leq n + 1$, and suppose that no box remains empty. Then either there is a box with precisely $N + 1$ items, or there is a set of boxes containing precisely N items.

Proof.

Divide the $N + n$ items, and suppose that no box remains empty. Delete one item from every box and consider the distribution of the remaining N items. Interpreting theorem B and “re-inserting” the deleted items in every box, gives the result. \square

3 Two sets of boxes.

Now consider the situation in which N items are divided over n red boxes and M items are divided over m green boxes. Is there a (non-trivial) set of red boxes that

contains exactly as many items as some set of green boxes. Clearly this is true for $N = M$, but in general the answer will depend on the way the items are actually distributed over the red and the green boxes respectively. Let a_i denote the number of items in the i^{th} red box ($1 \leq i \leq n$) and let b_j denote the number of items in the j^{th} green box ($1 \leq j \leq m$). The question raised can be rephrased as the decision problem for the predicate $EQ(a_1, \dots, a_n; b_1, \dots, b_m)$ which is defined to hold if and only if there is a set of a_i 's with exactly the same sum as some set of b_j 's (taking every "index" no more than once). We will also consider a "balanced" version of the problem, in the form of the predicate $BalEQ(a_1, \dots, a_n; b_1, \dots, b_m)$ which is defined to hold if and only if there is a set of k red boxes which together contain exactly as many items as some set of k green boxes (for some k with $1 \leq k \leq \min\{n, m\}$).

A desired solution to the problems raised would consist of some necessary and sufficient criteria for $EQ(x_1, \dots, x_n; y_1, \dots, y_m)$ and $BalEQ(x_1, \dots, x_n; y_1, \dots, y_m)$ which are easy to evaluate. In computational terms it means that we want an efficient algorithm for deciding arbitrary instances of the predicates EQ and $BalEQ$. Let NP denote the class of predicates that are decidable in nondeterministic polynomial time (cf. Garey & Johnson [1]). A problem or predicate E is called NP-complete whenever $E \in NP$ and every problem in NP can be polynomially transformed to E . (Thus NP-complete problems are the "hardest" problems in NP .) Many important and practical problems in computing have been shown to be NP-complete. As yet no deterministic polynomial time (i.e., "efficient") algorithm has been found for any of these problems. (For more information about the theory of NP-completeness we refer to [1].)

It is straightforward to verify that $EQ \in NP$ and $BalEQ \in NP$. For example, a nondeterministic algorithm for deciding $EQ(a_1, \dots, a_n; b_1, \dots, b_m)$ would proceed as follows: write down a choice of the a_i 's and a choice of the b_j 's, compute their sums, and output "true" if and only if the resulting values are equal. The executions of the algorithm are polynomial time bounded (in the size of the input, measured in bits). The argument for $BalEQ$ is similar. It is also straightforward to see that EQ is NP-complete, because it encompasses the so-called SUBSET-SUM problem ([1], problem SP13) which is already NP-complete. In fact we have the following result.

Theorem D.

- (i) EQ is NP-complete.
- (ii) $BalEQ$ is NP-complete.

Proof.

- (i) SUBSET-SUM is the following problem (cf. [1]): given a sequence of non-negative integers a_1, \dots, a_n and an integer b , is there a subsequence of the a_i 's that sums to exactly b . Clearly SUBSET-SUM is just equal to $EQ(x_1, \dots, x_n; y)$ and, because SUBSET-SUM is NP-complete, so is EQ .

(ii) *BalEQ* appears to be different from any known problem in [1], but it can be proved NP-complete as follows. Clearly $BalEQ \in NP$. There also is a polynomial transformation from SUBSET-SUM to *BalEQ*, as follows. Let a_1, \dots, a_n, b be an arbitrary instance of SUBSET-SUM. Let t be such that $2^t > \sum_1^n a_i \geq 2^{t-1}$. (Thus t is polynomial in the size of the problem, in bits.) Define $c_i = 2^{t+2i-1}$ and $d_i = e_i = 2^{t+2i-2}$ for $1 \leq i \leq n$, and consider the instance $BalEQ(a_1, \dots, a_n, c_1, \dots, c_n; b, d_1, \dots, d_n, e_1, \dots, e_n)$. We have to show that *BalEQ* holds in this instance if and only if the instance of SUBSET-SUM is solvable. Suppose *BalEQ* holds, i.e., there is a choice of k elements from $\{a_1, \dots, a_n, c_1, \dots, c_n\}$ and a choice of k elements from $\{b, d_1, \dots, d_n, e_1, \dots, e_n\}$ with exactly the same sum. If the choice from $\{a_1, \dots, a_n, c_1, \dots, c_n\}$ contains l of the c_i 's (for some $0 \leq l \leq k$) then the choice from $\{b, d_1, \dots, d_n, e_1, \dots, e_n\}$ necessarily contains l of the d_i 's and l of the e_i 's (thus a total of $2l$ elements). Thus we have $l < k$ and the choice from $\{a_1, \dots, a_n, c_1, \dots, c_n\}$ must contain some a_i 's. By inspecting the size of the numbers involved, the a_i 's that occur in the choice must sum to b . Conversely, suppose there is a solution to the instance of SUBSET-SUM. Let it consist of a choice of l of the a_i 's that sum to b , for some $1 \leq l \leq n$. Now the instance of *BalEQ* holds for $k = 2l - 1$: choose the same a_i 's and c_1, \dots, c_{l-1} from $\{a_1, \dots, a_n, c_1, \dots, c_n\}$, and choose b and d_1, \dots, d_{l-1} and e_1, \dots, e_{l-1} from $\{b, d_1, \dots, d_n, e_1, \dots, e_n\}$. The two sets of $2l - 1$ elements clearly have the same sum. The instance of *BalEQ* is obtained by a polynomial time computable transformation from the instance of SUBSET-SUM. It follows that *BalEQ* is NP-complete.

□

We now return to the more direct context of the pigeon hole principle. In a way the NP-completeness of *EQ* and *BalEQ* results from the “unbounded” size of the numbers involved, as we shall see. (This relates to the discussion of polynomial versus pseudo-polynomial time algorithms, cf. [1].) To get around it, we will now make the strict assumption that $1 \leq a_i \leq n$ ($1 \leq i \leq n$) and $1 \leq b_j \leq n$ ($1 \leq j \leq m$). We also assume that $n = m$. Let EQ' and $BalEQ'$ denote the correspondingly restricted versions of *EQ* and *BalEQ*.

Theorem E.

- (i) EQ' always holds.
- (ii) $BalEQ'$ is decidable by a polynomial time algorithm.

Proof.

- (i) We will prove the following, slightly stronger result: given $2n$ integers a_1, \dots, a_n and b_1, \dots, b_n in the range from 1 to n , there are $i, j \geq 1$ and $k, l \geq 0$ such that $a_i + \dots + a_{i+k} = b_j + \dots + b_{j+l}$. The proof involves a nice application of the pigeon hole principle, as follows.

Without loss of generality we may assume that $\sum_1^n a_i \leq \sum_1^n b_j$. For every i ($1 \leq i \leq n$) let j_i be the largest $j \geq 0$ such that $a_1 + \dots + a_i \geq b_1 + \dots + b_j$. (The j_i are well-defined.) Write $a_1 + \dots + a_i = b_1 + \dots + b_{j_i} + r_i$ and observe that $0 \leq r_i < n$ for every i . If $r_i = 0$ for some i , then we are done. Otherwise we have $1 \leq r_i < n$ for every i . As we now have n values in the range from 1 to $n - 1$, the pigeon hole principle can be used to assert that there must be indices i_1 and i_2 with $1 \leq i_1 < i_2 < n$ such that $r_{i_1} = r_{i_2}$. Clearly $j_{i_1} < j_{i_2}$, and in particular $j_{i_2} > 0$. It follows that $a_{i_1+1} + \dots + a_{i_2} = b_{j_1+1} + \dots + b_{j_2}$.

- (ii) To solve *BalEQ'* we have to determine whether there exists a choice of k elements from $\{a_1, \dots, a_n\}$ and a choice of k elements from $\{b_1, \dots, b_n\}$ with exactly the same sum, for some $k \geq 1$. (The problem is trivial when the two sets are not disjoint.) While the number of subsets is "large", it should be observed that every subset sum is bounded by n^2 . For $1 \leq t \leq n$ let $A_t(l, r)$ be defined to hold when there is some choice of l elements from $\{a_1, \dots, a_t\}$ that sums to r (for $0 \leq l \leq t$ and $0 \leq r \leq nt$). Similarly define $B_t(l, r)$ for the set $\{b_1, \dots, b_t\}$. Clearly A_1 and B_1 are trivial to compute. For $t \geq 1$, A_{t+1} can be computed from A_t as follows: if $r < a_t$ then $A_{t+1}(l, r) = A_t(l, r)$ otherwise $A_{t+1}(l, r) = A_t(l, r) \vee A_{t+1}(l-1, r-a_t)$. For consistency we let $A_t(0, 0) =$ "true" for all t . It easily follows that A_{t+1} can be computed from A_t in linear time and hence, A_n can be computed in polynomial time (for all arguments) from A_1 . Similarly B_n can be computed in polynomial time. Now *BalEQ'* holds if and only if there exist k and r with $1 \leq k \leq n$ and $1 \leq r \leq n^2$ such that both $A_n(k, r)$ and $B_n(k, r)$ hold. This is determined by a simple check of all possible values for k and r , in polynomial time. \square

Theorem E (i) shows that when N items are divided over n red boxes ($N \leq n^2$) and M items are divided over n green boxes ($M \leq n^2$) and all boxes get a number of items in the range from 1 to n , then there always is a (non-trivial) set of red boxes that together contain exactly as many items as some set of green boxes. It follows from the proof that the boxes in each set can in fact be chosen to be consecutive. Another interpretation is the following. Let $A = \{a_1, \dots, a_n\}$ and $B = \{b_1, \dots, b_n\}$ be two (possibly disjoint) sets of items with n items each, and let each item x have a weight $w(x)$ chosen from the integers ranging between 1 and n . Then there are non-trivial subsets $A' \subseteq A$ and $B' \subseteq B$ such that the elements from A' and B' can be exchanged without changing the total weight, i.e., $w(A) = w((A - A') \cup B')$ and $w(B) = w((B - B') \cup A')$, where $w(A) = \sum_{a \in A} w(a)$.

References

- [1] Garey, M.R., and D.S. Johnson, *Computers and Intractability: a Guide to the Theory of NP-Completeness*, W.H. Freeman & Co., San Francisco, 1979.
- [2] Hardy, G.H., and E.M. Wright, *An introduction to the Theory of Numbers*, 5th Edition, Oxford University Press, Oxford, 1979.
- [3] Liu, C.L., *Elements of Discrete Mathematics*, 2nd Edition, McGraw-Hill Book Comp., New York, 1985.