

# The $k$ -Equal Problem

---

MARTIN AIGNER

Freie Universität Berlin, Fachbereich Mathematik und Informatik  
Arnimallee 14, D-14195 Berlin  
(e-mail: aigner@math.fu-berlin.de)

*Received 30 December 2002; revised 16 June 2003*

Suppose we are given  $n$  coloured balls and an integer  $k$  between 2 and  $n$ . How many colour-comparisons  $Q(n, k)$  are needed to decide whether  $k$  balls have the same colour? The corresponding problem when there is an (unknown) linear order with repetitions on the balls was solved asymptotically by Björner, Lovász and Yao, the complexity being  $\theta(n \log \frac{2n}{k})$ . Here we give the exact answer for  $k > \frac{n}{2}$ :  $Q(n, k) = 2n - k - 1$ , and the order of magnitude for arbitrary  $k$ :  $Q(n, k) = \theta(\frac{n^2}{k})$ .

## 1. Introduction

Suppose we are given  $n$  coloured balls and an integer  $k$ ,  $2 \leq k \leq n$ . Two players, Paul and Carole, play the following game. At any stage Paul chooses two balls and asks whether they have the same colour, whereupon Carole answers ‘yes’ or ‘no’. The game ends when Paul either produces  $k$  balls coloured alike or states that no  $k$  balls have the same colour. How many questions  $Q(n, k)$  are needed in the worst case?

Two variants of this problem have been studied in detail. Suppose the  $n$  balls are coloured with two colours, and Paul has to produce a majority ball (or state that there is no majority). This variant is known as the majority problem and was solved first by Saks and Werman [9], and later by Alonso, Reingold, Schott [4, 5] and Wiener [10] using different methods. They answer it exactly:  $n - B(n)$  comparisons are needed, where  $B(n)$  is the number of 1s in the binary representation of  $n$ . When there is an arbitrary number of colours the solution  $\lceil \frac{3n}{2} \rceil - 2$  was given by Fisher and Salzberg [8]. A general account of this and related results is contained in Aigner [2, 3].

The other variant was considered by Björner, Lovász and Yao. Suppose we are given  $n$  real numbers (or equivalently a linear order on the balls with repetitions). How many

comparisons are necessary to decide whether  $k$  of the numbers are equal? They give the asymptotic answer  $\theta(n \log \frac{2^n}{k})$ .

In our situation we are at the other end. The coloured balls may be thought of as an antichain with repetitions. The general problem would then consist of an arbitrary poset (with repetitions). For example, for a linear order the case  $k = 2$  amounts to sorting the  $n$  numbers, while for an antichain we clearly have  $Q(n, 2) = \binom{n}{2}$ . For a general poset  $P$ , the case  $k = 2$  amounts to producing the poset  $P$ ; see, e.g., Aigner [1, Chapter 5].

Let us return to the topic of this note. It is to be expected that the cases  $k > \frac{n}{2}$  and  $k \leq \frac{n}{2}$  are of an entirely different nature. In the first case, a  $k$ -set of equally coloured balls is unique, while for  $k \leq \frac{n}{2}$  uniqueness is no longer guaranteed. This is reflected in the following two results.

**Theorem 1.1.** *Let  $k > \frac{n}{2}$ . Then  $Q(n, k) = 2n - k - 1$ .*

**Theorem 1.2.** *For general  $k$ ,  $2 \leq k \leq n$ , we have  $Q(n, k) = \theta(\frac{n^2}{k})$ , more precisely,  $\frac{1}{4} \frac{n^2}{k} < Q(n, k) < 2 \frac{n^2}{k}$ .*

**2. Proof of Theorem 1.1**

We use an argument similar to that in [8]. Let us look at the upper bound  $Q(n, k) \leq 2n - k - 1$  first. Paul uses the following algorithm.

**Phase 1.** Arrange the balls  $B_1, B_2, \dots, B_n$  in linear order, and compare the balls  $B_i$  one after the other. We set up a dynamic list  $L$  and a reservoir  $R$ . Initially,  $L = \{B_1\}$ ,  $R = \emptyset$ . Suppose that before  $B_j$  is compared the list is  $L = C_1 C_2 \dots C_s$  with reservoir  $R$ . Now compare  $B_j$  to the last ball  $C_s$  of the list. If they have the same colour, put  $B_j$  into  $R$ . Otherwise, enlarge  $L$  by moving  $B_j$  to the end and putting a ball  $Z$  of  $R$  behind  $B_j$  (in case  $R \neq \emptyset$ ). Phase 1 ends when the last ball  $B_n$  has been handled accordingly, so we have made  $n - 1$  comparisons so far.

The following facts are immediate.

- (a) All balls in  $R$  have the same colour which is equal to the colour of the last ball of  $L$ .
- (b) Neighbouring balls in  $L$  have different colours.
- (c) Since  $k > \frac{n}{2}$ , the only possible  $k$ -set of equally coloured balls contains the last ball of  $L$ .

**Phase 2.** Let  $L = D_t D_{t-1} \dots D_1 B A$  be the list with  $|R| = r$ , and thus  $t = n - 2 - r$ . Let us call the colour of  $A$  the *majority colour*, and all balls of this colour *majority balls*. Any ball with a colour different from  $A$  is called a *non-majority ball*, and we know that  $B$  is such a ball. If at the outset of Phase 2, either  $r + 1 \geq k$  or  $n + r + 2 - 2k \leq 0$ , then we are done. In the first case, we have found  $k$  majority balls, and in the second case there is no such set since the number of majority balls is at most

$$r + \left\lceil \frac{n-r}{2} \right\rceil = \left\lceil \frac{n+r}{2} \right\rceil \leq \frac{2k-2}{2} = k-1.$$

Now we scan the list  $L$  from the end, first comparing  $D_1$  with  $A$ , then  $D_2$  with  $A$ , and so on. Note that whenever  $D_i$  is a majority ball ( $i \geq 1$ ), then  $D_{i+1}$  is a non-majority ball. Suppose that up to a certain stage we have received  $\ell$  answers ‘yes’ and  $m$  answers ‘no’. The game is over when, for the first time, either after a ‘yes’ answer,

$$\ell + r + 1 \geq k, \tag{2.1}$$

or after a ‘no’ answer,

$$m + \ell + 1 + \left\lfloor \frac{n - r - 2 - 2\ell - m}{2} \right\rfloor \geq n - k + 1, \tag{2.2}$$

where  $\ell$  accounts for non-majority balls after each ‘yes’ answer, 1 for  $B$ , and the last summand for the remaining balls. Inequality (2.2) is equivalent to

$$m \geq n + r - 2k + 2. \tag{2.3}$$

It remains to estimate the number of questions  $\ell + m$ . Suppose (2.1) occurs first. Then  $\ell + r + 1 = k$  and  $2\ell - 1 + m \leq n - r - 2$  (note that the  $\ell$ th ‘yes’ may occur when  $D_\ell$  is compared to  $A$ ). Hence we obtain

$$\ell + m \leq n - r - 1 - \ell = n - k.$$

Suppose that (2.3) occurs first. Then  $\ell \leq k - 2 - r$  by (2.1) and  $m = n + r - 2k + 2$  by (2.3). Adding, we obtain again

$$\ell + m \leq n - k.$$

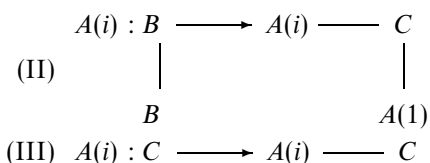
Hence Phase 2 uses at most  $n - k$  questions, and so  $Q(n, k) \leq 2n - k - 1$ .

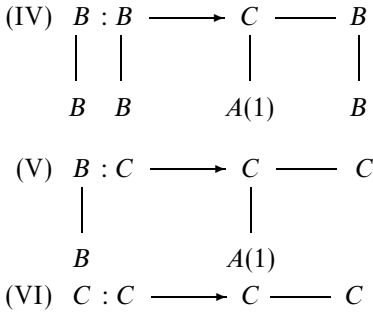
We come to the lower bound  $Q(n, k) \geq 2n - k - 1$ . Carole sets up the following data structure. She constructs a dynamic graph  $G$  on the balls where the edges correspond to the ‘no’ answers. The balls are dynamically labelled  $A, B$  or  $C$  with weights  $w$ , where all balls labelled  $B$  have weight  $1/2$  and all balls labelled  $C$  have weight  $0$ . Let  $\mathcal{A}, \mathcal{B}, \mathcal{C}$  be the sets of balls with labels  $A, B$  and  $C$ , respectively, and set  $\sum = \sum_{X \in \mathcal{A} \cup \mathcal{B}} w(X)$ . At the start all balls are labelled  $A$  with weight  $1$ , thus  $\sum = n$ . If  $X \in \mathcal{A}$  has weight  $i$ , then we write  $X = A(i)$ .

Now Carole answers according to the following rules, where we keep in mind that the  $A$ -balls aim for the  $k$ -majority.

- (I')  $A(i) : A(j)$  ‘yes’ if  $\sum = k$ , and the two balls are merged into a single ball  $A(i + j)$ .  
       ‘no’ if  $\sum > k$ , and an edge is inserted. In this case
- (I'')  $i = j = 1$  and we make the label change  $A(1) : A(1) \longrightarrow B - B$ .

In all other cases the answer is ‘no’, with the following label changes:





We note the following easy facts.

- (a) An edge of  $B$ -balls is only created when  $A(1) : A(1)$  are compared.
- (b) The  $B$ -balls constitute a matching, and there is no edge between different  $B$ -edges, between  $A(i)$  and  $A(j)$ , and between  $A(i)$  and  $B$ .
- (c)  $\Sigma$  is only reduced (by 1) in case (I'') and remains unchanged otherwise;  $\Sigma$  never drops below  $k$ . Furthermore, as long as  $\Sigma > k$ , all balls in  $\mathcal{A}$  have weight 1. We have  $|\mathcal{A}| \geq 1$  at any time, since  $\sum_{Y \in \mathcal{B}} w(Y) \leq \frac{n}{2} < k$ .
- (d) At a certain stage the number of ‘yes’ answers is  $\sum_{X \in \mathcal{A}} w(X) - |\mathcal{A}|$ . This follows from (I'). The number of ‘no’ answers (= edges) is at least  $2|\mathcal{C}| + \sum_{Y \in \mathcal{B}} w(Y)$ . This follows from (I'') to (VI).
- (e) At any stage a  $k$ -majority is possible. Indeed, since there are no edges between  $A : A$ ,  $A : B$  and different  $B$ -edges, all weighted balls in  $\mathcal{A}$  and one from each edge in  $\mathcal{B}$  may be coloured alike, summing to  $\Sigma \geq k$ .

**Claim.** As long as  $|\mathcal{A}| + |\mathcal{B}| \geq 2$ , there is also a non-majority possible, and hence the game is not over.

**Case 1:**  $\Sigma > k$ . Then all balls in  $\mathcal{A}$  have weight 1, and so all  $n$  balls may have different colours.

**Case 2:**  $\Sigma = k$ . If  $|\mathcal{B}| \geq 2$ , then  $\sum_{X \in \mathcal{A}} w(X) < k$ , and thus a non-majority is possible by assigning different colours to different  $A$ -balls and to all other balls. Suppose then  $\mathcal{B} = \emptyset$  and  $|\mathcal{A}| \geq 2$ . But here the same argument works since  $w(X) < k$  for all  $X \in \mathcal{A}$ .

Hence we conclude that the game is only over when  $|\mathcal{A}| = 1$ ,  $\mathcal{B} = \emptyset$ ,  $\Sigma = k$ , and thus  $|\mathcal{C}| = n - k$ . Invoking (d) we find that Paul has asked at least

$$k - 1 + 2(n - k) = 2n - k - 1$$

questions, and the proof is complete.

**Remarks.** (1) When  $k$  drops below  $\frac{n}{2}$ , the situation changes completely. Let  $n$  be even; then we have just seen that  $Q(n, \frac{n}{2} + 1) = \frac{3n}{2} - 2$ . By an argument similar to that of the proof of the theorem, it can be shown that  $Q(n, \frac{n}{2}) = 2n - 3$  for  $n \geq 6$ .

(2) In a variant which seems to favour Paul, Carole is only allowed  $c$  colours (unknown to Paul). It can be shown that in the strategy of Carole used in Theorem 1.1 she can force Paul to ask  $2n - k - 1$  questions using no more than 5 colours, for any  $n$ .

### 3. Proof of Theorem 1.2

We may restrict ourselves to  $k \leq \frac{n}{2}$ . Here the lower bound is easier. Carole always answers ‘no’. Let  $G$  be the dynamic graph on the  $n$  balls with edges corresponding to the comparisons. A  $k$ -set of identically coloured balls is possible as long as the independence number  $\alpha(G)$  is at least  $k$ . By Turán’s Theorem the minimal number of edges such that  $\alpha(G) \leq k - 1$  is achieved when  $G$  is the union of  $k - 1$  complete subgraphs of nearly equal size. Hence

$$Q(n, k) \geq (k - 1) \binom{\frac{n}{k-1}}{2} = \frac{1}{2} \left( \frac{n^2}{k-1} - n \right) > \frac{1}{4} \frac{n^2}{k},$$

where the last inequality follows from  $2k \leq n$ .

Now to the upper bound. Let  $n = qk + r$ ,  $0 \leq r < k$ . Paul arranges the balls  $1, 2, \dots, n$  around a circle in clockwise fashion, and proceeds in  $k - 1$  rounds. In the first round he compares  $i$  with  $i + 1, \dots, i + q \pmod n$  for all  $i$ . This takes  $qn$  comparisons. The ‘yes’ answers partition  $\{1, \dots, n\}$  into sets  $A_1, \dots, A_m$ , where balls in the same set have the same colour. Let us set  $\mathcal{C}(1) = \{A_1, \dots, A_m\}$ , with  $\sum_{i=1}^m |A_i| = n$ . Any set  $A \in \mathcal{C}(1)$  is an ‘arc’  $A = \{a_1, \dots, a_s\}$  read clockwise with  $a_{i+1} - a_i \leq q$ . Now either  $s \geq k$ , in which case we are finished, or the  $q$  balls following  $a_s$  have different colours from  $A$ . Furthermore, any ball  $b \notin A$  in the arc, that is,  $a_1 < b < a_s$ , also has a different colour. We say that  $A$  dominates at least  $q$  balls. Suppose that all sets  $A \in \mathcal{C}(1)$  have size  $\leq k - 1$ .

**Claim.** If the sets  $A_{i_1}, \dots, A_{i_r}$  in  $\mathcal{C}(1)$  form a possible  $\geq k$ -set with equal colour, then at least one of the sets contains at least two elements.

Suppose otherwise with  $A_{i_j} = \{a_j\}$ ,  $j = 1, \dots, k$ . Then by what we have just seen, any two consecutive balls are separated by at least  $q$  balls around the circle. It follows that  $n \geq kq + k$ , which contradicts  $r < k$ .

So, for the next round we need only compare sets  $A \in \mathcal{C}(1)$  with  $|A| \geq 2$  in clockwise fashion. In every round we compare  $A$  with the next  $q$  undecided balls. Let us fix some notation. After round  $i$  ( $1 \leq i \leq k - 1$ ),  $\mathcal{C}(i) = \{A_1, \dots, A_m\}$  is the partition into equally coloured sets determined by the ‘yes’ answers. We denote by  $s(A)$  the smallest number such that  $A \in \mathcal{C}(i)$  dominates at least  $s(A)q$  balls of different colour;  $\#A$  denotes the number of comparisons which involved balls of  $A$  in rounds 2 to  $i$ .

Now either there is  $A \in \mathcal{C}(i)$  with  $|A| \geq k$ , in which case we are finished, or we assume inductively that the following holds:

- (a)  $|A| \leq i \implies s(A) \geq |A|$ ,  $\#A \leq (|A| - 1)q$ ,
- (b)  $|A| \geq i + 1 \implies s(A) \geq i$ ,  $\#A \leq [\min(|A| - 2, s(A) - 1)]q$ .

Note that (a) and (b) hold for  $i = 1$ .

**Claim.** If, after round  $i$ , the sets  $A_1, \dots, A_\ell$  of  $\mathcal{C}(i)$  constitute a possible  $\geq k$ -set of equally coloured balls, then at least one of them contains at least  $i + 1$  balls.

Suppose not, and set  $f_t = \#\{j : |A_j| = t\}$ ,  $1 \leq t \leq i$ . Then  $\sum_{t=1}^i t f_t = |\bigcup_{j=1}^{\ell} A_j| \geq k$ . Invoking (a) we see that a  $t$ -set  $A$  is separated by at least  $tq$  balls from the next set, whence

$$n \geq \sum_{t=1}^i t f_t + \sum_{t=1}^i (t f_t) q \geq kq + k,$$

which cannot be.

Hence Paul need only compare sets  $A \in \mathcal{C}(i)$  with  $|A| \geq i + 1$  in round  $i + 1$ , and in fact he only compares all sets  $A \in \mathcal{C}(i)$  with

$$|A| \geq i + 1 \text{ and } s(A) = i.$$

We have to check the conditions (a) and (b) after round  $i + 1$ . Clearly, any set  $B \in \mathcal{C}(i + 1)$  is a union of sets  $A_1, \dots, A_h$  of  $\mathcal{C}(i)$ .

**Case 1:**  $h = 1$ , that is,  $B = A \in \mathcal{C}(i)$ . If  $|A| \leq i$  or  $|A| \geq i + 1$  and  $s(A) \geq i + 1$ , then  $A$  has not been compared,  $s(B) = s(A)$ , and (a) and (b) are satisfied. If  $|A| = i + 1$ ,  $s(A) = i$ , then  $A$  has been compared (with all  $q$  balls following  $A$  of different colour), hence  $s(B) = s(A) + 1$ . Thus  $s(B) = i + 1 = |B|$  and  $\#B = \#A + q \leq (|A| - 2)q + q = (|B| - 1)q$ , and (a) holds. Finally, if  $|A| \geq i + 2$ ,  $s(A) = i$ , then  $s(B) = s(A) + 1 = i + 1$ . Furthermore,  $\#A \leq [\min(|A| - 2, s(A) - 1)q] = (i - 1)q$ , thus  $\#B = \#A + q \leq iq \leq [\min(|B| - 2, s(B) - 1)q]$ , since  $|B| - 2 \geq i$  and  $s(B) - 1 = i$ .

**Case 2:**  $B = A_1 \cup \dots \cup A_h$  in clockwise order,  $h \geq 2$ . Since only  $A \in \mathcal{C}(i)$  with  $|A| \geq i + 1$ ,  $s(A) = i$  are compared, we have

$$\begin{aligned} |A_j| \geq i + 1, s(A_j) = i \text{ for } j = 1, \dots, h - 1, \\ |A_h| \geq i + 1, s(A_h) \geq i \text{ or } |A_h| \leq i, s(A_h) \geq |A_h|. \end{aligned}$$

It follows that

$$s(B) \geq (h - 1)i + \min(i, |A_h|) \geq i + 1.$$

If  $A_h$  is also compared, that is,  $|A_h| \geq i + 1$ ,  $s(A_h) = i$ , then

$$\begin{aligned} \#B &= \sum_{j=1}^h \#A_j + hq \leq \sum_{j=1}^h [\min(|A_j| - 2, i - 1)q] + hq \\ &= (i - 1)hq + hq = ihq \\ &\leq [\min(|B| - 2, s(B) - 1)]q, \end{aligned}$$

since  $|B| \geq (i + 1)h$ , and hence  $|B| - 2 \geq (i + 1)h - 2 \geq ih$  because of  $h \geq 2$ , and  $s(B) \geq i(h - 1) + (i + 1) = ih + 1$ , and thus  $s(B) - 1 \geq ih$ .

Finally, if  $A_h$  is not compared, then either  $|A_h| \geq i + 1$ ,  $s(A_h) \geq i + 1$  or  $|A_h| \leq i$ ,  $s(A_h) \geq |A_h|$ . This gives

$$\begin{aligned} \#B &= \sum_{j=1}^h \#A_j + (h - 1)q \leq \sum_{j=1}^{h-1} [\min(|A_j| - 2, i - 1)]q \\ &\quad + \begin{cases} [\min(|A_h| - 2, s(A_h) - 1)]q \\ (|A_h| - 1)q \end{cases} + (h - 1)q \\ &= i(h - 1)q + \begin{cases} [\min(|A_h| - 2, s(A_h) - 1)]q \\ (|A_h| - 1)q \end{cases} \\ &\leq [\min(|B| - 2, s(B) - 1)]q, \end{aligned}$$

since  $|B| \geq (i + 1)(h - 1) + |A_h|$ . Hence, with  $h \geq 2$ ,

$$|B| - 2 \geq i(h - 1) + h - 3 + |A_h| \geq i(h - 1) + |A_h| - 1,$$

and

$$s(B) \geq i(h - 1) + \begin{cases} s(A_h), \\ |A_h|, \end{cases}$$

and thus

$$s(B) - 1 \geq i(h - 1) + \begin{cases} s(A_h) - 1, \\ |A_h| - 1. \end{cases}$$

Performing rounds 2 to  $k - 1$ , we see that either we obtain a  $\geq k$ -set of equal colours, or that at the latest after round  $k - 1$  no such set is possible. Thus the game is finished after round  $k - 1$ . Counting the number of comparisons in rounds 2 to  $k - 1$  we find that after the last round  $\ell \leq k - 1$

$$\sum_{A \in \mathcal{C}(\ell)} \#A \leq \sum_{A \in \mathcal{C}(\ell)} (|A| - 1)q < nq.$$

Altogether, this gives, with round 1, at most  $2nq \leq 2\frac{n^2}{k}$  comparisons, and the proof is complete. □

### References

- [1] Aigner, M. (1988) *Combinatorial Search*, Wiley.
- [2] Aigner, M. (2004) Variants of the majority problem. *Discrete Math.* **137** 3–25.
- [3] Aigner, M. Two colors and more. Preprint.
- [4] Alonso, L., Reingold, E. and Schott, R. (1993) Determining the majority. *Inform. Process. Lett.* **47** 253–255.
- [5] Alonso, L., Reingold, E. and Schott, R. (1997) Average-case complexity of determining the majority. *SIAM J. Comput.* **26** 1–14.

- [6] Björner, A., Lovász, L. and Yao, A. (1992) Linear decision trees: Volume estimates and topological bounds. In *Proc. 24th ACM Symp. on Theory of Computing*, ACM Press, New York, pp. 170–177,
- [7] Du, D. and Hwang, F. (1993) *Combinatorial Group Testing*, World Scientific.
- [8] Fisher, M. and Salzberg, S. (1982) Finding a majority among  $n$  votes. *J. Algorithms* **3** 375–379.
- [9] Saks, M. and Werman, M. (1991) On computing majority by comparisons. *Combinatorica* **11** 383–387.
- [10] Wiener, G. (2002) Search for a majority element. *J. Statist. Plann. Inference* **100** 313–318.