

# On Penrose's square-root law and beyond

Werner Kirsch

ABSTRACT. In certain bodies, like the Council of the EU, the member states have a voting weight which depends on the population of the respective state. IN this article we ask the question which voting weight guarantees a 'fair' representation of the citizens in the union. The traditional answer, the square-root law by Penrose, is that the weight of a state (more precisely: the voting power) should be proportional to the square-root of the population of this state. The square root law is based on the assumption that the voters in every state cast their vote independently of each other. In this paper we concentrate on cases where the independence assumption is not valid.

## 1. Introduction

All modern democracies rely on the idea of representation. A certain body of representatives, a parliament for example, makes decisions on behalf of the voters. In most parliaments each of its members represents roughly the same number of people, namely the voters in his or her constituency.

There are other bodies in which the members represent different numbers of voters. A prominent example is the Council of the European Union. Here ministers of the member states represent the population of their respective country. The number of people represented in the different states differs from about 400,000 for Malta to more than 82 million for Germany. Due to this fact the members of the Council have a certain number of votes depending on the size of the country they represent, e.g. 3 votes for Malta, 29 votes for Germany. The votes of a country cannot be split, but have to be cast as a block.<sup>1</sup>

Similar voting systems occur in various other systems, for example in the Bundesrat, Germany's state chamber of parliament and in the electoral college in the USA.<sup>2</sup>

---

<sup>1</sup>The current voting system in the Council is based on the treaty of Nice. It has additional components to the procedure described above, which are irrelevant in the present context. For a description of this voting system and further references see e.g. [5].

<sup>2</sup>The electoral college is not exactly a heterogeneous voting system in the sense defined below, but it is very close to it.

Let us call such a system in which the members represent subsystems (states) of different size a *heterogeneous voting system*. In the following we will call the assembly of representatives in a heterogeneous voting system the *council*, the sets of voters represented by the council members the *states*.

It is quite clear, that in a heterogeneous voting system a bigger state (by population) should have at least as many votes in the council as a smaller state. It may already be debatable whether the bigger states should have *strictly* more votes than the smaller states (cf. the Senate in the US constitution). And if yes, how much more votes the bigger state should get?

In this note we address the question: ‘What is a fair distribution of power in a heterogeneous voting system?’

There exist various answers to this question, depending on the interpretation of the words ‘fair’ and ‘power’.

The usual and quite reasonable way to formulate the question in an exact way is to use the concept of power indices. One calls a heterogeneous voting system fair if all voters in the member states have the same influence on decisions of the council. By ‘same influence’ we mean that the power index of each voter is the same regardless of her or his home state. If we choose then Banzhaf power index to measure the influence of a voter we obtain the celebrated Penrose’s square-root law (see e.g. [3]).

The square-root law states that the distribution of power in a heterogeneous voting system is fair if the power (index) of each council member  $i$  is proportional to  $\sqrt{N_i}$ , where  $N_i$  is the population of the state which  $i$  represents.

In their book [3] Felsenthal and Machover formulate a second square-root law. There they base the notion of ‘fairness’ on the concept of *majority deficit*.

The majority deficit is zero if the voters favoring the decision of the council are the *majority*. If the voters favoring the decision of the council are the *minority* then the majority deficit is the margin between the number of voters objecting to the decision and those agreeing with it (see Def. 3.3.16 in [3]).

The notion of fairness we propose in this paper is closely related to the concept of majority deficit. We will call a decision of the council *in agreement with the popular vote* if the percentage of voters agreeing with a proposal (popular vote) is as close as possible to the percentage of council votes in favor of the proposal. (We will make this notion precise in the next section.)

For both concepts we have to average over the possible voting configurations. This is usually done by assuming that voters vote independently of each other. The main purpose of this note is to investigate some (we believe reasonable) models where voters do not vote independently.

We will discuss two voting models with voting behavior which is *not* independent. The first model considers societies which have some kind of ‘common belief’. A typical situation of this kind is a strong religious group

(or church) influencing the voting behavior of the voters. This model is discussed in detail in Section 3.

In the other model voters tend to vote the same way ‘the majority does’. This is a situation where voters do not want to be different from others. We call this the *mean field model* referring to an analogous model from statistical physics. See Section 5 for this model.

In fact, both models can be interpreted in terms of statistical physics. Statistical physics considers (among many other things) magnetic systems. The elementary magnet, called a spin, has two possible states which are ‘+1’ or ‘-1’ (spin up, spin down). This models voting ‘yes’ or ‘no’ in a voting system. Physicists consider different kinds of interactions between the single spins, one given through an exterior magnetic-field - corresponding to a society with ‘a common belief’ - or through the tendency of the spins to align - corresponding to the second voting model. We discuss the analogy of voting models with spin systems in Section 4.

Our investigations of voting models with statistical dependence is much inspired by the paper [7]. The first model is also based on the work by Straffin [8].

It does not come as a surprise that we obtain a square-root law for a model with independent voters, just as in the case considered by Felsenthal and Machover ([3]).

For the mean field model we still get a square-root law for the best possible representation in the council *as long as* the mutual interaction between voters is not too strong.

However as the coupling between voters exceeds a certain threshold, the fairest representation in the council is no longer given by votes proportional to  $\sqrt{N_i}$  but rather by votes proportional to  $N_i$ . This is a typical example of a phase transition.

In the model of common belief the fair representation weight depends on the strength of the common belief for large populations. If this strength is independent of the population size fair representation is almost always given by voting weights proportional to  $N_i$ , the square-root law occurring only in marginal cases. However, if the common belief decreases with increasing population one can get any power law behavior  $N_i^\alpha$  for the optimal weight as long as  $\frac{1}{2} \leq \alpha \leq 1$ . In fact, statistical investigations on real life data suggest that this might happen (see [4]).

We leave the mathematical proofs of our results for the appendices (Sections 7 to 9).

**Acknowledgment:** It is a pleasure to thank Hans-Jürgen Sommers, Duisburg-Essen and Wojciech Słomczyński and Karol Życzkowski, Krakow for valuable discussions.

## 2. The general model

We consider  $N$  voters, denoted by  $1, 2, \dots, N$ . Each of them may vote ‘yes’ or ‘no’; abstentions are not allowed. The vote of the voter  $i$  is denoted by  $X_i$ .

The possible voting results are  $X_i = +1$  representing ‘yes’ and  $X_i = -1$  for ‘no’. We consider the quantity  $X_i$  as random, more precisely there is a probability measure  $\mathbb{P}$  on the space  $\{-1, 1\}^N$  of possible voting results. This measure will be specified later. The conventional assumption on  $\mathbb{P}$  is that the random quantities  $X_i$  are independent from each other, but we are *not* making this assumption here.

Our interpretation of this model is as follows. The voters react on a proposal in a rational way, that is to say: A voter does *not* roll a dice to determine his or her voting behavior but he or she votes for or against a given proposal according to his/her personal belief, knowledge, experience etc. It is rather the proposal which is the source of randomness in this system. We imagine the voting system is fed with propositions in a completely random way. This could be either a real source of proposals or just a Gedankenexperiment to measure the behavior of the voting system.

The rationality of the voters implies that a voter who casts a ‘yes’ on a certain proposition will necessarily vote ‘no’ on the diametrically opposed proposition. Since we assume that the proposals are completely random any proposal and its antithetic proposal must have the same probability. This implies

$$\mathbb{P}(X_i = 1) = \mathbb{P}(X_i = -1) = \frac{1}{2} . \quad (2.1)$$

More generally, we conclude that

$$\mathbb{P}(X_{i_1} = \xi_1, \dots, X_{i_r} = \xi_r) = \mathbb{P}(X_{i_1} = -\xi_1, \dots, X_{i_r} = -\xi_r) \quad (2.2)$$

for any set  $i_1, \dots, i_r$  of voters and any  $\xi_1, \dots, \xi_r \in \{-1, 1\}$ .

We call the property (2.2) the *symmetry* of the voting system. Any measure  $\mathbb{P}$  satisfying (2.2) is called a *voting measure*.

The symmetry assumption (2.2) does *not* fix the probability measure  $\mathbb{P}$ . Only if we assume in addition that the  $X_i$  are statistically independent we can conclude from (2.2) that

$$\mathbb{P}(X_{i_1} = \xi_1, \dots, X_{i_r} = \xi_r) = \left(\frac{1}{2}\right)^r . \quad (2.3)$$

So far, we have not specified any decision rule for the voting system. The above probabilistic setup is completely independent from the voting rule, a fact which was emphasized in the work [7].

A *simple majority rule* for  $X_1, \dots, X_N$  is given by the decision rule: Accept a proposal if  $\sum_{j=1}^N X_j > 0$  and reject it otherwise.

By a *qualified majority rule* we mean that at least a percentage  $q$  (called the *quota*) of votes is required for the acceptance of a proposal. In term of the  $X_j$  this means:

$$\sum_{j=1}^N X_j \geq (2q - 1)N. \quad (2.4)$$

Indeed, it is not hard to see that the number of affirmative votes is given by

$$\frac{1}{2} \left( \sum_{j=1}^N X_j + N \right).$$

From this the assertion (2.4) follows.

In particular, the simple majority rule is obtained form (2.4) by choosing  $q$  slightly bigger than  $\frac{1}{2}$ .

The sum  $\sum_{j=1}^N X_j$  gives the difference between the number of ‘yes’-votes and the number of ‘no’-votes. We call the quantity

$$M(X) := \left| \sum_{j=1}^N X_j \right| \quad (2.5)$$

the margin of the voting outcome  $X = (X_1, \dots, X_N)$ . It measures the size of the majority with which the proposal is either accepted or rejected in simple majority voting.

In qualified majority voting with quota  $q$  the corresponding quantity is the  $q$ -margin  $M_q(X)$  given by:

$$M_q(X) := \left| \sum_{j=1}^N X_j - (2q - 1)N \right|. \quad (2.6)$$

Now, we turn to voting in the council. We consider  $M$  states, the state number  $\nu$  having  $N_\nu$  voters. Consequently the total number of voters is  $N = \sum N_\nu$ . The vote of the voter  $i$  in state  $\nu$  is denoted by  $X_{\nu i}$ ,  $\nu = 1, \dots, M$  and  $i = 1, \dots, N_\nu$ .<sup>3</sup>

We suppose that each state government knows the opinion of (the majority of) the voters in that state and acts accordingly.<sup>4</sup> That is to say: If the majority of people in state  $\nu$  supports a proposal, i.e. if

---

<sup>3</sup>We label the states using Greek characters and the voters within a state by Roman characters.

<sup>4</sup>Although this is the central idea of representative democracy this idealization may be a little naive in practice.

$$\sum_{i=1}^{N_\nu} X_{\nu i} > 0 \quad (2.7)$$

then the representative of state  $\nu$  will vote ‘yes’ in the council otherwise he or she will vote ‘no’. If we set  $\chi(x) = 1$  for  $x > 0$ ,  $\chi(x) = -1$  for  $x \leq 0$  the representative of state  $\nu$  will vote

$$\xi_\nu = \chi \left( \sum_{i=1}^{N_\nu} X_{\nu i} \right) \quad (2.8)$$

in the council. If the state  $\nu$  has got a weight  $w_\nu$  in the council the result of voting in the council is given by:

$$\sum_{\nu=1}^M w_\nu \xi_\nu = \sum_{\nu=1}^M w_\nu \chi \left( \sum_{i=1}^{N_\nu} X_{\nu i} \right). \quad (2.9)$$

Thus, the council’s decision is affirmative if  $\sum_{\nu=1}^M w_\nu \xi_\nu$  is positive, provided the council votes according to simple majority rule.

The result of a popular vote in all countries  $\nu = 1, \dots, N$  is

$$P = \sum_{\nu=1}^M \sum_{i=1}^{N_\nu} X_{\nu i}. \quad (2.10)$$

We will call voting weights  $w_\nu$  for the council *fair* or *optimal*, if the council’s vote is as close as possible to the public votes. To make this precise let us define

$$C = \sum_{\nu=1}^M w_\nu \chi \left( \sum_{i=1}^{N_\nu} X_{\nu i} \right) \quad (2.11)$$

the result of the voting in the council. Both  $P$  and  $C$  are random quantities which depend on the random variables  $X_{\nu i}$ . So, we may consider the mean square distance  $\Delta$  between  $P$  and  $C$ , i.e. denoting the expectation over the random quantities by  $\mathbb{E}$ , we have

$$\Delta = \mathbb{E}((P - C)^2) \quad (2.12)$$

$$= \mathbb{E} \left( \left\{ \sum_{\nu=1}^M \sum_{i=1}^{N_\nu} X_{\nu i} - \sum_{\nu=1}^M w_\nu \chi \left( \sum_{i=1}^{N_\nu} X_{\nu i} \right) \right\}^2 \right). \quad (2.13)$$

In a democratic system the decision of the council should be as close as possible to the popular vote, hence we call a system of weights *fair* or *optimal* if  $\Delta = \Delta(w_1, \dots, w_M)$  is minimal among all possible values of  $w_\nu$ .

In the following we suppose that the random variables  $X_{\nu i}$  and  $X_{\mu j}$  are independent for  $\nu \neq \mu$ . This means that voters in different states are not

correlated. We do not assume at the moment that two voters from the same state vote independently of each other.

We have the following result:

THEOREM 2.1. *Fair voting in the council is obtained for the values*

$$\begin{aligned} w_\nu &= \frac{1}{2} \mathbb{E} \left( \left| \sum_{i=1}^{N_\nu} X_{\nu i} \right| \right) \\ &= \frac{1}{2} \mathbb{E} \left( M(X_{\nu 1}, \dots, X_{\nu N_\nu}) \right). \end{aligned}$$

This result can be viewed as an extension of Penrose's square-root law to the situation of correlated voters. We will see below that it gives  $w_\nu \sim \sqrt{N_\nu}$  for independent voters.

Theorem 2.1 has a very easy - we hope convincing - interpretation:  $w_\nu$  is the expected margin of the voting result in state  $\nu$ . In other words, it gives the expected number of people in state  $\nu$  that agree with the voting of  $\nu$  in their council minus those that disagree, i.e. the net number of voters which the council member of  $\nu$  actually represents.

If we choose any multiple  $cw_1, \dots, cw_{N_\nu}$  ( $c > 0$ ) of the weights  $w_1, \dots, w_{N_\nu}$  we obtain the same voting system as the one defined by  $w_1, \dots, w_n$ . In this sense the weight  $w_\nu$  of Theorem 2.1 are not unique, but the voting system is.

We will prove Theorem 2.1 in section 7. We remark that the proof requires the symmetry assumption (2.2) and the independence of voters from *different* states.

The next step is to compute the expected margin  $\mathbb{E}(M(X_\nu))$ , at least asymptotically for large number of voters  $N_\nu$ . This quantity depends on the correlation structure between the voters in state  $\nu$ . As we will see, different correlations between voters give very different results for  $\mathbb{E}(M(X_\nu))$  and hence for the optimal weight  $w_\nu$ .

We begin with the classical case of independent voters.

THEOREM 2.2. *If the voters in state  $\nu$  cast their votes independently of each other then*

$$\mathbb{E} \left( \left| \sum_{i=1}^{N_\nu} X_{\nu i} \right| \right) \sim c \sqrt{N_\nu} \tag{2.14}$$

for large  $N_\nu$ .

Thus, we recover the square-root law as we expected. (For the square-root law see Felsenthal and Machover [3].) In terms of power indices the independence assumption is associated to the Banzhaf power index. Therefore, it is not surprising that also the Banzhaf index leads to a square-root rule.

It is questionable (as we know from the work of Gelman, Katz and Bafumi [4]) whether the independent voters model is valid in many real-life voting systems. This is one of the reasons to extend the model as we do in the present paper.

### 3. The ‘common belief’ model

In this section we consider a model we dub the ‘common belief model’. It generalizes a voting measure introduced and investigated by Straffin [8] in connection with the Shapley-Shubik power index.

We imagine that inside a certain society there is a strong common belief which may for example be due to a powerful religious group, a generally accepted political ideology or a strong tradition. This causes a tendency to a creation of strong majorities in a certain type of questions. For example, in a country with a strong catholic majority there may be a strongly correlated view about abortion among voters, but, may be, not about speed limits on highways. One might have a similar effect if a person dominates the public or private media or both.

We model such a situation by introducing a random variable  $Z$  which reflects the ‘common belief’ on the subject at hand.  $Z = +1$  means that all voters agree to accept the given proposal,  $Z = -1$  means that all voters will reject it. The random variable  $Z$  is allowed to take any value in  $[-1, 1]$ . If  $Z = 0$  there is no common belief on the proposal. If  $Z > 0$  there is some common belief favoring the proposal which is weak if  $Z$  is close to 0 and strong if it is close to 1. The probability distribution of  $Z$  is denoted by  $\mu$ , hence

$$\mu([a, b]) = \mathbb{P}(Z \in [a, b]). \quad (3.1)$$

$Z$  has to satisfy a symmetry condition similar to (2.1), namely

$$\mathbb{P}(Z \in [a, b]) = \mathbb{P}(Z \in [-b, -a]), \quad (3.2)$$

i.e.

$$\mu([a, b]) = \mu([-b, -a]) . \quad (3.3)$$

In our model the ‘common belief’ variable  $Z$  influences the probability that a voter  $i$  votes  $\pm 1$ . Given  $Z = \zeta \in [-1, 1]$ , then the conditional probability that  $X_i = 1$  given  $Z = \zeta$  is given by:

$$\mathbb{P}(X_i = 1 | Z = \zeta) = \frac{1}{2}(\zeta + 1) = p_\zeta . \quad (3.4)$$

Thus, if  $\zeta = 1$  any voter  $i$  will vote ”+1” with probability 1, if  $\zeta = -1$  he or she will vote ”+1” with probability 0 and if  $\zeta = 0$  then  $i$  votes with probability one half in favor or against the proposal.

We denote by  $P_s$  the probability measure on  $\{-1, +1\}^N$  with  $P_s(X_i = 1) = s$  and  $P_s(X_i = -1) = 1 - s$ , which makes the  $X_i$  independent. We also denote by  $E_s$  the expectation with respect to  $P_s$ . Note that the probability  $p_\zeta = \frac{1}{2}(1 + \zeta)$  in (3.4) is chosen in such a way that  $E_{p_\zeta}(X_i) = \zeta$ . Thus the

value of the ‘common belief’ variable  $Z$  gives the expected voting result of a single voter.

For *given*  $Z = \zeta$  we assume the  $X_i$  to be independent. Thus we have

$$\mathbb{P}(X_1 = \xi_1, \dots, X_N = \xi_N) = \int \left( \prod_1^N P_{p_\zeta}(X_i = \xi_i) \right) d\mu(\zeta) . \quad (3.5)$$

The measure  $\mathbb{P}$  in (3.5) depends on the probability distribution  $\mu$ , hence we sometimes denote it by  $\mathbb{P}_\mu$ .

The probability measure  $\mathbb{P}_\mu$  defined in (3.5) satisfies the symmetry condition (2.2) due to assumption (3.3). Of course,  $\mathbb{P}_\mu$  defines a whole class of examples, each (symmetric) probability measure  $\mu$  on  $[-1, 1]$  defines its unique  $\mathbb{P}_\mu$ . If we choose  $\mu = \delta_0$ , i.e.  $\mu([a, b]) = 1$  if  $a \leq 0 \leq b$  and  $= 0$  otherwise, we obtain independent random variables  $X_i$  as discussed in the final part of section 2. Indeed,  $\mu = \delta_0$  means that  $Z = 0$ , consequently (3.5) defines independent random variables. Observe, that this is the only measure for which  $Z$  assumes a fixed value ( $\mu$  has to be symmetric!).

Another interesting example is the case when  $\mu$  is the uniform distribution on  $[-1, 1]$ . This case was considered by Straffin [8]. He observed that this model is intimately connected with the Shapley-Shubik power index.

To apply the ‘common belief’ model to a given heterogeneous voting model we have to specify the measure  $\mu$ , of course. In fact, this measure may change from state to state. In particular, one may argue that larger states tend to have a less homogeneous population and hence a weaker influence of a religious or political group. For example, we will later discuss a model modifying Straffin’s example where  $\mu(dz) = \frac{1}{2}\chi_{[-1,1]}(z)dz$  to a measure where  $\mu_N$  depends on the population  $N$ , namely

$$\mu_N(dz) = \frac{1}{2a_N}\chi_{[-a_N, a_N]}(z)dz \quad (3.6)$$

with parameters  $0 < a_N \leq 1$ . In particular, if we have  $a_N \rightarrow 0$  as  $N \rightarrow \infty$ , the parameter  $a_N$  reflects the tendency of a common belief to decrease with a growing population.

Except for the trivial case  $\mu = \delta_0$  the random variables  $X_i$  are never independent under  $\mathbb{P}_\mu$ . This can be seen from the covariance

$$\langle X_i, X_j \rangle_\mu := \mathbb{E}_\mu(X_i X_j) - \mathbb{E}_\mu(X_i)\mathbb{E}_\mu(X_j) . \quad (3.7)$$

In (3.7) as well as in the following  $\mathbb{E}_\mu$  denotes expectation with respect to  $\mathbb{P}_\mu$ . In fact, the random variables  $X_i$  are always positively correlated:

**THEOREM 3.1.** *For  $i \neq j$  we have*

$$\langle X_i, X_j \rangle_\mu = \int \zeta^2 d\mu(\zeta) . \quad (3.8)$$

The quantity  $\int \zeta^2 d\mu(\zeta)$  is called the second moment of the measure  $\mu$ . Since the first moment  $\int \zeta d\mu(\zeta)$  vanishes due to (3.3) the second moment equals the variance of  $\mu$ . Observe that  $\int \zeta^2 d\mu(\zeta) = 0$  implies  $\mu = \delta_0$ . For

independent random variables  $\langle X_i, X_j \rangle_\mu = 0$ , so (3.8) implies that  $X_i, X_j$  depend on each other unless  $\mu = \delta_0$ .

To investigate the impact of the common belief measure  $\mu$  on the ideal weight in a heterogeneous voting model we have to compute the quantity

$$\mathbb{E}_\mu(|\sum X_i|) \quad (3.9)$$

for a measure  $\mu$  and population  $N$  (at least for large  $N$ ). This is done with the help of the following Theorem:

$$\text{THEOREM 3.2. } \left| \mathbb{E}_\mu\left(\frac{1}{N} \left| \sum_1^N X_i \right| \right) - \int |\zeta| d\mu(\zeta) \right| \leq \frac{1}{\sqrt{N}}.$$

If we choose  $\mu \neq \delta_0$  independent of the (population of the) state Theorem 3.2 implies that the optimal weight in the council is *proportional* to  $N$  (rather than  $\sqrt{N}$ ). This is true in particular for the original Straffin model [8] where  $\mu_n \equiv \frac{1}{2} \chi_{[-1,1]}(z) dz$  which corresponds to the Shapley-Shubik power index.

Let us define  $\bar{\mu} = \int |\zeta| d\mu(\zeta)$ . If  $\mu = \mu_N$  depends on the population then

$$\mathbb{E}_{\mu_N}(|\sum X_i|) \sim N \bar{\mu}_N$$

as long as  $\bar{\mu}_N \geq \frac{1}{N^{1/2-\varepsilon}}$  for some  $\varepsilon > 0$ . However, if  $\bar{\mu}_N \leq \frac{1}{N^{1/2-\varepsilon}}$ , then

$$E_{\mu_N}(|\sum X_i|) \sim \sqrt{N}.$$

Hence, in this case we rediscover a square-root law.

We summarize:

**THEOREM 3.3.** *Let us suppose that a state with a population of size  $N$  is characterized by a common belief measure  $\mu_N$ , then:*

(1) *If*

$$\bar{\mu}_N = \int |\zeta| d\mu_N(\zeta) \geq C \frac{1}{N^{1/2-\varepsilon}} \quad (3.10)$$

*for some  $\varepsilon > 0$  and for all large  $N$  then the optimal weight  $w_N$  is given by:*

$$w_N = \mathbb{E}_\mu(|\sum_1^N X_i|) \sim N \bar{\mu}_N. \quad (3.11)$$

(2) *If*

$$\bar{\mu}_N = \int |\zeta| d\mu_N(\zeta) \leq C \frac{1}{N^{1/2+\varepsilon}} \quad (3.12)$$

*then for large  $N$  the optimal weight  $w_N$  is given by:*

$$w_N = \mathbb{E}_\mu(|\sum_1^N X_i|) \sim \sqrt{N}. \quad (3.13)$$

**Example:** In our Straffin-type example (3.6) we choose:

$$\mu_N(dz) = \frac{1}{2a_N} \chi_{[-a_N, a_N]}(z) dz, \quad (3.14)$$

then:

$$\bar{\mu}_N = \frac{1}{2} a_N. \quad (3.15)$$

So, if  $a_N \leq C \frac{1}{\sqrt{N}}$  we have  $w_N \sim \sqrt{N}$ , otherwise we obtain  $w_N \sim a_N$ .

REMARKS 3.4.

- (1) *Our result shows that in all cases the optimal weight  $w_N$  satisfies  $C\sqrt{N} \leq w_N \leq N$ . It is a matter of empirical studies to determine which measure  $\mu_N$  is appropriate to the given voting system. Any of the empirical results of [4] can be modeled by an appropriate choice of  $\mu_N$ .*
- (2) *It is only  $\bar{\mu}_N$  that enters the formulae (3.11) and (3.13), no other information about  $\mu_N$  is relevant. The quantities  $\bar{\mu}_N$  can be estimated using Theorem 3.2. In fact, more is true by the following result.*

THEOREM 3.5. *Let  $P_N$  be the distribution of  $\frac{1}{N} \sum_{i=1}^N X_i$  under the measure  $\mathbb{P}_{\mu_N}$  then the sequence of measures  $P_N - \mu_N$  converges weakly to 0.*

Note that the distribution of  $\frac{1}{N} \sum_{i=1}^N X_i$  is the distribution of the voting results of the voter  $i = 1, \dots, N$ . This is the quantity considered in [4]. Theorem 3.5 tells us that the distribution of the voting results for large number  $N$  of voters is approximately equal to the distribution  $\mu_N$ . In particular, for independent voting the voting result is always extremely tight while for Straffin's example any voting result has the same probability, i.e. it is equally likely that a proposal gets 99% or 53% of the votes.

#### 4. Voting models as spin systems

Spin systems are a central topic in statistical physics. They model magnetic phenomena. The spin variables, usually denoted by  $\sigma_i$ , may take values in the set  $\{-1, +1\}$  with  $+1$  and  $-1$  meaning 'spin up' and 'spin down' respectively. The spin variables model the elementary magnets of the material (say the electrons or nuclei in a solid). The index  $i$  runs over an index set  $I$  which represents the set of elementary magnets.

The probability measure underlying the statistical structure is typically given by a 'Gibbs measure' defined through an energy functional  $\mathcal{E}(\{\sigma_i\}_{i \in I})$ .  $\mathcal{E}$  gives the energy of a given spin configuration  $\{\sigma_i\}$ . The system prefers configurations with low energy  $\mathcal{E}$ . This is expressed in the Gibbs measure given by:

$$q(\{\sigma_i\}_{i \in I}) = e^{-\beta \mathcal{E}(\{\sigma_i\}_{i \in I})}. \quad (4.1)$$

The parameter  $\beta$  plays the role of an inverse temperature.  $q$  defines a (counting) measure on the space  $\Omega = \{-1, +1\}^I$ . It has total mass:

$$\mathcal{Z} = \sum_{\{\sigma_i\} \in \Omega} e^{-\beta \mathcal{E}(\{\sigma_i\}_{i \in I})}. \quad (4.2)$$

Hence we obtain a *probability* measure by setting:

$$p(\{\sigma_i\}_{i \in I}) = \mathcal{Z}^{-1} e^{-\beta \mathcal{E}(\{\sigma_i\}_{i \in I})}. \quad (4.3)$$

Of course, we may interpret any spin system as a voting system with voting measure  $p$ , as long as  $\mathcal{E}(\{\sigma_i\}) = \mathcal{E}(\{-\sigma_i\})$ , and vice versa.

In particular, independent voting corresponds to the energy functional  $\mathcal{E}(\{\sigma_i\}) \equiv 1$ .

Moreover, the ‘common belief’ model is given by an energy function:

$$\mathcal{E}(\{\sigma_i\}) = -h \sum_i \sigma_i \quad (4.4)$$

where  $h$  is a random variable connected to the variable  $Z$  defined in (3.1) by:

$$\frac{1}{2}(1 + Z) = \frac{e^h}{e^h + e^{-h}}. \quad (4.5)$$

Note, that when  $h$  runs from  $-\infty$  to  $\infty$  in (4.5) the value of  $Z$  runs monotonously from  $-1$  to  $+1$ .

In term of statistical physics the above model is a system without spin-spin interaction in a random but constant magnetic field. The inverse temperature  $\beta$  which we encountered in equation (4.1) is superfluous in this model as it can be absorbed in the magnetic field strength  $h$ .

## 5. The voters’ interaction model

In the common belief model the voting behavior of each voter is influenced by a preassigned, a priori given common belief variable  $Z$ . The correlation between the voters results from the general voting tendency described by the value of  $Z$ .

In this section we investigate a model with a *direct* interaction between the voters, namely a tendency of the voters to vote in agreement with each other. In the view of statistical physics this corresponds to the tendency of magnets to align. There are various models in statistical physics to prescribe such a situation. Presumably the best known one is the Ising model where neighboring spins interact in the prescribed ways. The neighborhood structure is most of the time given by a lattice (e.g.  $\mathbb{Z}^d$ ). The results on the system depend strongly on that neighborhood structure, in the case of the lattice  $\mathbb{Z}^d$  on the dimension  $d$ .

In the following we consider another, in fact easier model where no such assumption on the local ‘neighborhood’ structure has to be made. We

consider it an advantage of the model that very little of the microscopic correlation structure of a specific voting system enters into the model.

The model we are going to consider is known in statistical mechanics as the *Curie-Weiss model* or the *mean field model* (see e.g. [9], [1] or [2]). In this model a given voter (spin) interacts with all the other voters (resp. spins) in a way which makes it more likely for the voters (spins) to agree than to disagree. This is expressed through an energy function  $\mathcal{E}$  which is smaller if voters agree. Note that a *small* energy for a given voting configuration (relative to the other configurations) leads to a *high* probability of that configuration relative to the others through formula (4.3).

The energy  $\mathcal{E}$  for a given voting outcome  $\{X_i\}_{i=1\dots N}$  is given in the mean field model by:

$$\mathcal{E}(\{X_i\}) = -\frac{J}{N-1} \sum_{\substack{i,j \\ i \neq j}} X_i X_j. \quad (5.1)$$

Here  $J$  is a non negative number called the coupling constant. According to (5.1) the energy contribution of a single voter  $X_i$  is expressed through the averaged voting result of all other voters  $\frac{1}{N-1} \sum_{j \neq i} X_j$ . If  $X_i$  agrees in sign with this average the voter  $i$  makes a negative contribution to the total energy, otherwise  $X_i$  will increase the total energy. The strength of this negative or positive contribution is governed by the coupling constant  $J$ . In other words: situations for which  $X_i$  agrees with the other voters in average are more likely than others. This can be seen from the formula for the probability of a given voting outcome, namely:

$$p_J(\{X_i\}) = \mathcal{Z}^{-1} e^{-\mathcal{E}(\{X_i\})} = \mathcal{Z}^{-1} e^{\frac{J}{N-1} \sum_{i \neq j} X_i X_j} \quad (5.2)$$

where we have set

$$\mathcal{Z} = \sum_{\{X_i\} \in \{\pm 1\}^N} e^{-\mathcal{E}(\{X_i\})}. \quad (5.3)$$

As before the parameter  $\beta$  is not needed, it can be absorbed in the coupling constant  $J$ .

Our goal is to compute the average:

$$w_N = \mathbb{E}_{J,N} \left( \left| \sum_{i=1}^N X_i \right| \right). \quad (5.4)$$

Here  $\mathbb{E}_{J,N}$  denotes expectation with respect to the measure defined in (5.2). The quantity  $w_N$  gives the optimal weight in the council for a population of  $N$  voters with a correlation structure given by a mean-field model with coupling constant  $J$ . We will see that the value of  $w_N$  changes dramatically when  $J$  changes from a value below one to a value above one. This has to do with the fact that the mean-field model undergoes a phase transition at the point  $J = 1$  (see [1, 2, 9]).

THEOREM 5.1.

(1) If  $J < 1$  then

$$w_N = \mathbb{E}_{J,N}(|\sum_{i=1}^N X_i|) \sim \frac{\sqrt{2}}{\sqrt{\pi}} \frac{1}{\sqrt{1-J}} \sqrt{N} \quad \text{as } N \rightarrow \infty. \quad (5.5)$$

(2) If  $J > 1$  then

$$w_N = \mathbb{E}_{J,N}(|\sum_{i=1}^N X_i|) \sim C(J) N \quad \text{as } N \rightarrow \infty. \quad (5.6)$$

REMARKS 5.2.

- (1) By  $x_N \sim y_N$  as  $N \rightarrow \infty$  we mean that  $\lim_{n \rightarrow \infty} \frac{x_N}{y_N} = 1$ .  
(2) The constant  $C(J)$  in (5.6) can be computed: If  $J > 1$  then  $C(J)$  is the (unique) positive solution  $C$  of

$$\tanh(JC) = C. \quad (5.7)$$

Note that for  $J \leq 1$  there is no positive solution of equation of 5.7.

The proof of Theorem 5.1 will be given in section 9.

## 6. Conclusions

The above calculations show that one can reproduce the square-root law as well as the results of [4] and other laws by assuming particular correlation structures among the voters of a certain country. To find the right model is a question of adjusting the parameters of the models to empirical data of the country under consideration. Moreover, the models allow us to investigate questions about voting systems on a theoretical level. We believe that the models described above can help to understand voting behavior in many situations.

To *design* a nonhomogeneous voting system for a *constitution* in the light of our results is a question of different nature. Even knowing the correlation structure of the countries in question exactly would be of limited value to design a constitution. Constitutions are meant for a long term period, correlation structures of countries on the other hand are changing even on the scale of a few years.

One might argue that modern societies have a tendency to decrease the correlation between their members. In all modern states, at least in the West, the influence of churches, parties, and unions is constantly declining.

In addition to this it seems more important to protect small countries against a domination of the big ones than the other way round. This motivates us to choose a square-root law in these long term cases.

## 7. Appendix 1: Proofs for section 2

We start with a short Lemma:

LEMMA 7.1. *Suppose  $X_1, \dots, X_N$  are  $\{-1, 1\}$ -valued random variables with the symmetry property (2.2) then*

$$\mathbb{E}\left(\sum_{i=1}^N X_i\right) = 0 \quad (7.1)$$

and

$$\mathbb{E}\left(\sum_{i=1}^N X_i \chi\left(\sum_{i=1}^N X_i\right)\right) = \frac{1}{2} \mathbb{E}\left(\left|\sum_{i=1}^N X_i\right|\right). \quad (7.2)$$

REMARK 7.2. *As defined above  $\chi(x) = 1$  if  $x > 0$ ,  $\chi(x) = -1$  if  $x \leq 0$ .*

PROOF. (2.2) implies

$$\mathbb{P}(X_i = 1) = \mathbb{P}(X_i = -1) = \frac{1}{2}$$

hence  $\mathbb{E}(X_i) = 0$  and (7.1) follows.

To prove (7.2) we observe that due to (2.2)

$$\begin{aligned} \mathbb{E}\left(\left|\sum_{i=1}^N X_i\right|\right) &= \mathbb{E}\left(\sum_{i=1}^N X_i \chi\left(\sum_{i=1}^N X_i\right)\right) - \mathbb{E}\left(\sum_{i=1}^N X_i \chi\left(-\sum_{i=1}^N X_i\right)\right) \\ &= 2\mathbb{E}\left(\sum_{i=1}^N X_i \chi\left(\sum_{i=1}^N X_i\right)\right). \end{aligned}$$

□

We turn to the proof of Theorem 2.1.

PROOF. (Theorem 2.1) Let us abbreviate:  $S_\nu := \sum_{i=1}^{M_\nu} X_{\nu i}$ .

Observe that the  $S_\nu$  are independent by assumption and satisfy  $\mathbb{E}(S_\nu) = 0$ , moreover

$$\mathbb{E}(S_\nu \chi(S_\mu)) = 0 \text{ if } \nu \neq \mu \quad (7.3)$$

and

$$\mathbb{E}(S_\nu \chi(S_\nu)) = \frac{1}{2} \mathbb{E}(|S_\nu|) \quad (7.4)$$

by Lemma 7.1. To find the minimum of the function

$$\Delta(w_1, \dots, w_M) = \mathbb{E}\left(\left(\sum_{\nu=1}^M S_\nu - \sum_{\nu=1}^M w_\nu \chi(S_\nu)\right)^2\right)$$

we look at the zeros of  $\frac{\partial \Delta}{\partial w_\mu}$ .

$$\begin{aligned} 0 = \frac{\partial \Delta}{\partial w_\mu} &= -2\mathbb{E}\left(\left(\sum_{\nu=1}^M S_\nu - \sum_{\nu=1}^M w_\nu \chi(S_\nu)\right) \chi(S_\mu)\right) \\ &= -2\mathbb{E}(S_\mu \chi(S_\mu) - w_\mu \chi(S_\mu) \chi(S_\mu)). \end{aligned}$$

So

$$w_\mu \mathbb{E}((\chi(S_\mu))^2) = \mathbb{E}(S_\mu \chi(S_\mu)) = \frac{1}{2} \mathbb{E}(|S_\mu|) .$$

Since  $\chi(S_\mu)^2 = 1$  we obtain

$$w_\mu = \frac{1}{2} \mathbb{E}(|S_\mu|) .$$

□

We turn to the proof of Theorem 2.2.

PROOF. Let  $X_1, \dots, X_N$  be  $\{-1, 1\}$ -valued random variables with  $P(X_i = 1) = P(X_i = -1) = \frac{1}{2}$ . Then

$$\mathbb{E}(|\sum_1^N X_i|) = \sqrt{N} \mathbb{E}(|\frac{1}{\sqrt{N}} \sum_1^N X_i|) .$$

By the central limit theorem (see e.g. [6])  $\frac{1}{\sqrt{N}} \sum_1^N X_i$  has asymptotically a normal distribution with mean zero and variance 1, hence

$$\mathbb{E}(|\frac{1}{\sqrt{N}} \sum_1^N X_i|) \rightarrow \frac{\sqrt{2}}{\sqrt{\pi}} .$$

□

### 8. Appendix 2: Proofs for Section 3

PROOF. (Theorem 3.1) Since  $\mathbb{E}_\mu(X_i) = 0$ ,

$$\begin{aligned} \langle X_i, X_j \rangle_\mu &= \mathbb{E}_\mu(X_i X_j) & (8.1) \\ &= \mathbb{P}_\mu(X_i = X_j = 1) + \mathbb{P}_\mu(X_i = X_j = -1) - 2\mathbb{P}_\mu(X_i = 1, X_j = -1) \\ &= \int d\mu(\zeta) \{ P_{\frac{1}{2}(1+\zeta)}(X_i = X_j = 1) + P_{\frac{1}{2}(1+\zeta)}(X_i = X_j = -1) \\ &\quad - 2P_{\frac{1}{2}(1+\zeta)}(X_i = 1, X_j = -1) \} \\ &= \int d\mu(\zeta) \{ \frac{1}{4}(1+\zeta)^2 + \frac{1}{4}(1-\zeta)^2 - \frac{1}{2}(1-\zeta^2) \} \\ &= \int \zeta^2 d\mu(\zeta) . \end{aligned}$$

□

To prove Theorem 3.2 we need the following Lemma:

LEMMA 8.1.  $\mathbb{E}_\mu(\frac{1}{N} |\sum (X_i - Z)|) \leq \frac{1}{\sqrt{N}}$ .

PROOF.

$$\begin{aligned} \mathbb{E}_\mu(\frac{1}{N} |\sum (X_i - Z)|) &= \frac{1}{N} \mathbb{E}_\mu(|\sum (X_i - Z)|) \\ &\leq \frac{1}{N} \left\{ \mathbb{E}_\mu\left(\left(\sum (X_i - Z)\right)^2\right) \right\}^{1/2} \\ &= \frac{1}{N} \left\{ \int d\mu(\zeta) E_{p_\zeta}\left(\left(\sum_1^N (X_i - \zeta)\right)^2\right) \right\}^{1/2} \quad (8.2) \end{aligned}$$

Given  $Z = \zeta$  the random variables  $X_i - \zeta$  have mean zero and are independent with respect to the measure  $P_{p_\zeta}$ , thus

$$\mathbb{E}_{p_\zeta} \left( \left( \sum_1^N (X_i - \zeta) \right)^2 \right) = N \mathbb{E}_{p_\zeta} (X_i - \zeta)^2 = N(1 - \zeta^2) \leq N,$$

hence

$$(8.2) \leq \frac{1}{\sqrt{N}} \left( \int d\mu(\zeta) (1 - \zeta^2) \right)^{1/2} \leq \frac{1}{\sqrt{N}}.$$

□

Using Lemma 8.1 we are in a position to prove Theorem 3.2:

PROOF. (1) Suppose that:

$$\bar{\mu}_N = \int |\zeta| d\mu_N(\zeta) \geq C \frac{1}{N^{1/2-\varepsilon}} \quad (8.3)$$

then we estimate:

$$\begin{aligned} \mathbb{E}_{\mu_N} \left( \frac{1}{N} \left| \sum_1^N X_i \right| \right) &= \mathbb{E}_{\mu_N} \left( \left| \frac{1}{N} \sum_1^N (X_i - Z) + Z \right| \right) \\ &\leq \mathbb{E}_{\mu_N} (|Z|) + \mathbb{E}_{\mu_N} \left( \left| \frac{1}{N} \sum_1^N (X_i - Z) \right| \right) \\ &\leq \bar{\mu}_N + \frac{1}{\sqrt{N}} \end{aligned} \quad (8.4)$$

by Lemma 8.1. Moreover

$$\begin{aligned} \mathbb{E}_{\mu_N} \left( \frac{1}{N} \left| \sum_1^N X_i \right| \right) &\geq \mathbb{E}_{\mu_N} (|Z|) - \mathbb{E}_{\mu_N} \left( \left| \frac{1}{N} \sum_1^N X_i - Z \right| \right) \\ &\geq \bar{\mu}_N - \frac{1}{\sqrt{N}}. \end{aligned} \quad (8.5)$$

Hence

$$\left| \mathbb{E}_{\mu_N} \left( \frac{1}{N} \left| \sum_1^N X_i \right| \right) - \bar{\mu}_N \right| \leq \frac{1}{\sqrt{N}} \quad (8.6)$$

which proves (3.11).

(2) To prove (3.13) we obtain by the same reasoning as above:

$$\left| \mathbb{E}_{\mu_N} \left( \frac{1}{N} \left| \sum_1^N X_i \right| \right) - \bar{\mu}_N \right| \leq \frac{1}{\sqrt{N}} \quad (8.7)$$

□

We end this section with the proof of Theorem 3.5:

PROOF. We have to prove that for bounded continuous functions  $f$ :

$$\int \left( f\left(\frac{1}{N} \sum_{i=1}^N X_i\right) - f(Z) \right) d\mathbb{P}_{\mu_N} \rightarrow 0. \quad (8.8)$$

The convergence (8.8) is clear for continuously differentiable  $f$  from Lemma 8.1. It follows for arbitrary bounded continuous  $f$  by a density argument.  $\square$

### 9. Appendix 3: Proofs for section 5

In this section we prove Theorem 5.1.

PROOF. (Theorem 5.1 (1))

We denote by  $E_0^{(N)}$  the expectation of the coin tossing model for  $N$  independent symmetric  $\{+1, -1\}$ -valued random variables, i.e.:

$$E_0^{(N)}(F(X_1, \dots, X_N)) = \frac{1}{2^N} \sum_{\{x_i\} \in \{+1, -1\}^N} f(x_1, \dots, x_N). \quad (9.1)$$

We set:

$$\mathcal{Z}_{JN} = E_0^{(N)} \left( e^{\frac{J}{2} \left( \frac{1}{\sqrt{N}} \sum_{i=1}^N X_i \right)^2} \right) \quad (9.2)$$

and:

$$\mathcal{X}_{JN} = E_0^{(N)} \left( \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N X_i \right| e^{\frac{J}{2} \left( \frac{1}{\sqrt{N}} \sum_{i=1}^N X_i \right)^2} \right). \quad (9.3)$$

Then:

$$\mathbb{E}_{JN} \left( \left| \sum_{i=1}^N X_i \right| \right) = \sqrt{N} \frac{\mathcal{X}_{JN}}{\mathcal{Z}_{JN}}. \quad (9.4)$$

Under the probability law  $E_0^{(N)}$  the random variables  $X_i$  are centered and independent, thus the central limit theorem (see e.g. [6]) tells us that  $\frac{1}{\sqrt{N}} \sum_{i=1}^N X_i$  converges in distribution to a standard normal distribution. Consequently, for  $J < 1$  and  $N \rightarrow \infty$ :

$$\mathcal{Z}_{JN} \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(1-J)x^2}{2}} dx = \frac{1}{\sqrt{1-J}} \quad (9.5)$$

and:

$$\mathcal{X}_{JN} \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} |x| e^{-\frac{(1-J)x^2}{2}} dx = \frac{\sqrt{2}}{\sqrt{\pi}} \frac{1}{1-J}. \quad (9.6)$$

Consequently:

$$\mathbb{E}_{JN} \left( \left| \sum_{i=1}^N X_i \right| \right) = \sqrt{N} \frac{\mathcal{X}_{JN}}{\mathcal{Z}_{JN}} \sim \frac{\sqrt{2}}{\sqrt{\pi}} \frac{1}{\sqrt{1-J}} \sqrt{N}. \quad (9.7)$$

$\square$

PROOF. (Theorem 5.1 (2))

By Theorem 6.3 in [1] the distribution  $\nu_N$  of  $S_N = \frac{1}{N} \sum_{i=1}^N X_i$  converges weakly to the measure  $\nu = \delta_{-C(J)} + \delta_{C(J)}$  where  $C(J)$  was defined in (5.7).

Hence,

$$\mathbb{E}_J\left(\left|\sum_{i=1}^N X_i\right|\right) = N \mathbb{E}_J(|S_N|) \quad (9.8)$$

$$= N \int |\lambda| d\nu_N(\lambda) \quad (9.9)$$

$$\approx N \int |\lambda| d\nu(\lambda) \quad (9.10)$$

$$= N C(J). \quad (9.11)$$

□

## References

- [1] Bolthausen, E.; Sznitman, A.: Ten lectures on random media, Birkhäuser (2002).
- [2] Dorlas, T.: Statistical Mechanics, Institute of Physics Publishing (1999).
- [3] Felsenthal, D.; Machover, M.: The measurement of power: theory and practice, problems and paradoxes, Edward Elgar (1998).
- [4] Gelman, A.; Katz, J.; Bafumi, J.: Standard voting power indexes don't work: an empirical analysis, British Journal of Political Science 34, 657–674 (2004).
- [5] Kirsch, W.: On the distribution of power in the Council of ministers of the EU. Preprint, available from <http://www.ruhr-uni-bochum.de/mathphys/>
- [6] Lamperti, J.: Probability, Wiley (1996).
- [7] Laruelle, A.; Valenciano, F.: Assessing success and decisiveness in voting situations, Social Choice and Welfare 24, 171-197 (2005).
- [8] Straffin, P.: Power indices in politics, Brams et. al. (Eds.): Political and related models, Springer (1982).
- [9] Thompson, C.: Mathematical Statistical Mechanics, Princeton University Press (1972).

INSTITUT FÜR MATHEMATIK, RUHR-UNIVERSITÄT BOCHUM,  
D-44780 BOCHUM, GERMANY

*E-mail address:* `werner.kirsch@ruhr-uni-bochum.de`