# Behaviour, Interaction and Control of User Communities

Matthew Collinson

**Abstract:**

Most modern organisations have information security policies that are designed to guide the behaviour of their user communities. It is often impractical for these policies to be enforced directly, and users frequently have incentives not to comply. In both realistic and simplified situations the resulting principal-agent problem can be extremely complicated. Consequently, managers often have to make decisions about security policy in the face of a high degree of uncertainty, both about user behaviour and the ambient threat environment. The purpose of this paper is to draw attention to some of the complexities using a variety of types of model, and to suggest ways in which progress towards practical, model-based decision processes might be made. No single model - or type of model - is likely to provide complete insight into the problem. First to be considered is a decision-making process using calculation of utility, and based on inferences about population behaviour derived from empirical data. The issues surrounding a practical methodology featuring simulation are discussed. The use of game theory is considered as a way of understanding the interaction between an organization and its users. It is further proposed that methods from statistical mechanics can be used to provide models of interaction and influence within the user community - these suggest that extreme non-linearities may be present in the behaviour of the community. In each case, attention is paid to the difficulties of collecting the data required by the models.

# Behaviour, Interaction and Control of User Communities

Matthew Collinson

collinson@hp.com

HP Labs, Long Down Avenue, Bristol, BS34 8QZ, UK

February 22, 2010

### Abstract

Most modern organisations have information security policies that are designed to guide the behaviour of their user communities. It is often impractical for these policies to be enforced directly, and users frequently have incentives not to comply. In both realistic and simplified situations the resulting principal-agent problem can be extremely complicated. Consequently, managers often have to make decisions about security policy in the face of a high degree of uncertainty, both about user behaviour and the ambient threat environment.

The purpose of this paper is to draw attention to some of the complexities using a variety of types of model, and to suggest ways in which progress towards practical, model-based decision processes might be made. No single model — or type of model — is likely to provide complete insight into the problem. First to be considered is a decision-making process using calculation of utility, and based on inferences about population behaviour derived from empirical data. The issues surrounding a practical methodology featuring simulation are disussed. The use of game theory is considered as a way of understanding the interaction between an organisation and its users. It is further proposed that methods from statistical mechanics can be used to provide models of interaction and influence within the user community — these suggest that extreme non-linearities may be present in the behaviour of the community. In each case, attention is paid to the difficulties of collecting the data required by the models.

## 1 Introduction

Information security management in most organisations involves the *control* [1] of a human population — the community of *users*. People are, of course, intelligent, autonomous agents with their own preferences and goals, and abilities to anticipate, and adapt to, the contol mechanisms of the organisation. Consequently, control of the user community is a major part of the security management problem.

Security management is itself part of a larger problem, namely that of achieving the goals of the organisation. An organisation may have to trade-off security against other goals, for example, maximization of revenue and minimization of costs. Its portfolio of investments in security also constitute a trade-off between security objectives, such as confidentiality, integrity and availability. It is by now widely recognised that many of the security decision problems faced by an organisation need to be addressed using the techniques of economics [3, 28, 29, 38].

Organisations may have a range of mechanisms through which they may attempt to exercise control. First, there are *technological mechanisms* that enforce or preclude certain forms of user behviour. For example, users normally have privileges that are restricted using an access control system. Second, there are *(security) policies* that are used to mandate or forbid certain behaviours. For example, users may be forbidden from using certain applications or from allowing others to use their account, but required to keep their password secret and to lock their client machine

---

[1]The word 'control' is used throughout this paper the sense of applied mathematics and engineering [15] and, in general, connotes effective influence rather than puppetry.

before leaving it out-of-sight. Third, there are *reward* and *punishment mechanisms* that can be used to encourage desirable, and discourage undesirable, user behaviour. A user may be fired or suspended for the abuse of systems. An organisation will typically use different combinations of these three types of mechanism to achieve different security objectives. For example, it may have a policy that all users must immediately apply emergency security patches to their clients, block access to any user who fails to comply, and use disciplinary measures against a user who tries to overcome the block.

In addition to the above mechanisms, an organisation may hope that individual ethics, social norms and innate human capacities help to constrain user behaviour in the desired way. Unfortunately, these forces may have the opposite effect. The social and behavioural sciences must be crucial to understanding and improving the effectiveness of security mechanisms [1, 2, 3, 6, 35, 70, 71].

Technological mechanisms are not always available to achieve security objectives, nor are they always desirable. Their use is, in effect, a *centralised* mechanism for control. By contrast, policy (often together with punishment) is a useful *decentralised* control: some aggregation of the users' autonomous actions determine the risk faced by the firm. For one thing, this may have benefits in terms of cost of implementation. For another, it may minimise the impact on other of the organisations goals. For example, users may be able to choose when to apply a security patch to their machines, minimising disruption and so loss of productivity. It has even been suggested that (in some situations) users could be allowed to choose their own access privileges [40]. On the other-hand, it allows users opportunities not to comply, and instead to engage in behaviour that exposes the organisation to additional risk. Indeed, users may have significant incentives not to comply, even where they are not malicious insiders. If the organisation does not have the ability to check-up on (and punish) users then it faces an example of a *principal-agent* problem.

In reality, the picture is much more complex. There may be a mix of centralised and decentralised control used to achieve security objectives: for example, if users do not patch their own systems before a deadline, it may be done automatically. Whilst disciplinary mechanisms exist in theory, in practice it may be that they are only used after extreme security events.

This paper is focussed on policy mechanisms, discussing the resulting behaviour of user communities and trying to tie this to its economic impact. This is done in an effort to provide an objective decision process for the organisation regarding the use of such mechanisms. The viewpoint taken is rather high-level, looking at different modelling methodologies for studying compliance in populations, rather than focussing on specific technologies and situations.

The fundamental problem studied throughout the paper is the same: how policy should be set in order that the resulting behaviour of the user community is of maximum benefit to the organisation. This is evidently a problem of control. What differs between sections is the level of detail in models and the methodologies applied. Given the complexity of the problem it is unlikely that a single methodology or model can give a comprehensive view.

In Section 2, utility-based decision-making for security policy is discussed using an example. The assumptions and data required for such a model are examined in some detail. The use of simulation as an augmentation of this method for practical (time and resource-constrained) decision-making is discussed.

The responses of intelligent users to policy are hard to predict. Section 3 focusses on the use of techniques from game theory to model decision-making processes in the context of rational, individual user responses guided by preferences.

A further source of complexity in community behaviour is the social influence of one user upon another. Section 4 is centred on the use of statistical mechanics to model the behaviour of a user community, paying special attention to the importance of dynamic interaction between users. In Section 5 models of community behaviour in the presence of anticipative interaction between individual users are discussed.

All models in science and engineering are only approximate pictures of the real world (perhaps excepting some models of fundamental physics). The models in this paper are extreme simplifications and mathematical idealisations of real situations. Nevertheless, very rich and complex structure remains. One must be very careful about inferring general conclusions from idealised models. The theoretical models herein are treated as only suggestive. Much more data is needed

2

to affirm their pragmatism. It is one thing to produce plausible mathematical accounts and post hoc rationalizations of natural phenomena, but quite another to be able to make useful predictions. The purpose of this essay is to draw attention to the complexities of the decision problem, to discuss possible modelling methodologies, and to suggest where interesting data and results may be found.

# 2 Simple decisions from quantitative data on populations

## 2.1 A utility-based view

Let us begin with a very simple stylised example in order to start the discussion of how security decisions about user communities might ideally be made. In line with much of economic theory this is based around the idea that the firm should optimize an objective function that captures its preferences.

Consider a situation in which a firm [2] considers the implementation of one of a range of new security policies designed to achieve a similar effect (this may include the option of having no policy). Assume a fixed number of users of size $N$. Users have the option to act in a way that is insecure (for the firm). For simplicity, let us use the word '*comply*' generically for the behaviour of users who do not have insecure behaviour (even when there is no policy for them to defy). Rather than considering the behaviour of individual agents, the firm cares only about the total number (or frequency) of compliant users. Further, suppose that the firm knows the response of the user community to each of the firm's policy choices. Thus there is a map taking each policy $x$ to a distribution over user behaviours. Since user choices are binary, the distribution can be encapsulated in a single figure $f(x) \in [0, 1]$ giving the frequency of compliant users. Further still, suppose that it is known that the risk to the firm is an increasing function $r(m)$ of the number $m$ of non-compliant users.

The preferences of the firm are described by a *loss function*. The loss function $\mathcal{L}_0(x)$ for the firm is a function of the policy $x$. Suppose that each policy $x$ has a cost $c(x)$ to the firm that includes both implementation costs and loss of productivity (perhaps caused by loss of productivity for users, or loss of system availabilty). The loss should be increasing in both the frequency $f(x)$ and the cost $c(x)$. The firm faces a trade-off between risk $r(Nf(x))$ and the costs $c(x)$. Given all of the above, the firm could, in principle, take a decision $x$ by minimizing its loss $\mathcal{L}_0(x)$. The details of the loss function are not required here, since no optimization will be performed in this paper.

It should be evident that to construct such simple example and such a simple decision process has required a very heavy burden of assumptions:

*Isolation of policy.* The policy decision is discussed above as though it is taken in isolation: the firm only cares about the implementation costs and productivity loss appearing in the loss function. In practice, it may be difficult to do this — policies overlap and user-responses to additional security burden are not expected to be additive. It is the marginal effect of new policy that will typically be of interest. The model must reflect this fact, or it must be otherwise justifiable that the effects can be treated in isolation.

*Time.* The time over which the model is applied has been condensed into a single, discrete point. Multi-period or continuous time models would often be more realistic.

User decisions have been treated as irreversible, and the times at which a user complies (or does not) are irrelevant to both the cost and risk calculations. If the decision is irreversible, an uptake curve depicting the frequency of compliance $f(x)(t)$ at each time $t$ is required. A situation resembling the standard model of diffusion of (technological) innovations [57] would not be unexpected here, with its well-known categories of innovators, early adopters, early majority, late majority, and laggards. See also [35] on possible uses for this theory in security management. For reversible or repeated user decisions (for example, daily decisions about encryption of data on

---

[2]The terminolgy 'firm' is preferred to 'organisation' for most of the remainder of the paper. This is solely for the sake of brevity — it is not required to be a commercial enterprise.

portable media) the uptake curve may not be monotonically increasing. Geographical and social ties will affect the spread of compliance. The loss function for the firm would be determined by the time-parametrised path followed by the user community (through the space of possible community configurations). The subject of cultural effects and social dynamics is explored in more detail in Section 4.

In more general models the threat environment may be able exploit both slow uptake and time periods when user behaviour varies (for example, holidays and corporate events).

*Population assumptions.* The risks and costs that the firm faces (in this model) are driven only by the number of non-compliant users. The structure of the population and the distribution of roles and privileges is irrelevant. In reality, different user groups have different privileges, are subject to different policies and exist in heterogeneous cultures with differing social norms. The risks and costs for the firm vary greatly from user to user.

The user population is assumed to be static. For a real firm there may be large variations in user population over the time period of interest. The rates of arrival and departure of users are likely to have a strong influence on population behaviour.

*Exogenised user behviour* Users are rational agents. However, in the model above their behaviour has been treated as exogenous, with the assumption that overall population reaction $f(x)$ to each policy choice $x$ is known. In addition, one should not overlook the possibility of insider threats from users. A more complex picture involving user preferences and their ability to anticipate may be required. A loss function for users may have to be complicated if it is to accurately represent preferences. This points towards the use of game-theoretic techniques. The treatment of user behaviour starting from their preferences is discussed in Section 3 below.

*Exogenised risk.* In the model, the risk to the firm is given as a known function $r$ of the number of compliant users. A firm with security that is perceived to be weak will likely be attacked more. The risk function would have to take into account how this happens: in particular, $r$ may grow very sharply at some threshold value of non-compliance.

In reality, the risk comes from a complex and rapidly-changing *threat environment*. The threat-environment is driven by rational attackers who interact with each other and with legitimate organisations. Changes in the threat-environment should be understood as part of a co-evolutionary process with the rest of the digital and social world.

The use of game-theoretic ideas may be helpful here. The firm must be able to value the assets it wishes to protect, both to itself and its attackers, and understand the mechanisms that prevent attackers from realising their goal. This has been studied in some detail by other authors, see for example [64, 34]. The fact that attackers often focus on the softest of a range of target firms means that any firm would in principle have to take those other firms' behaviours into account. For example, a bank with (perceived) poor security arrangements may be unable to compete for customers for online-accounts. The threat couples each firm's decision to that of many others. In principle, a model may have to incorporate all of these other agents' potential actions and preferences.

The simple decision process outlined above will not usually be achievable in practice because either the assumptions cannot be met, or data is required that cannot feasibly be collected.

Simple data anticipating the behaviour of attackers that consitute the threat environment quickly becomes out-of-date. Data about other firms' security policies is confidential and inaccessible. Even if it weren't, one would need to anticipate attackers behaviour in order to know which data to collect from other firms. Consequently, the form for the risk function would be very hard to generate. The dynamics of social interaction and cultural effects can make population behaviour very hard to describe, so the function $f$ used to give the population's reaction may be difficult to estimate.

Suppose that the difficulties above could be at least partially overcome, and that one could express the preferences of the firm using a loss of the form $\mathcal{L}_0(x_0, m)$, where $x_0$ is the policy choice and $m$ is the frequency of compliant users. If for each $x_0$ the population reaction $f(x_0)$ could be

given by a probability distribution, then the firm may choose to calculate the expected loss

$$\mathbf{E}(\mathcal{L}_0(x_0)) = \sum_{m \in \mathbb{N}} \mathcal{L}_0(x_0, m).\, Pr(N.f(x_0) = m \mid x_0)$$

and make its decision by minimizing this quantity. More generally, where the structure of compliance in the population matters, the firm may wish to consider its expected loss

$$\mathbf{E}(\mathcal{L}_0(x_0)) = \sum_{x_1, \ldots, x_N} \mathcal{L}_0(x_0, x_1, \ldots, x_N).\, Pr(x_1, \ldots, x_N \mid x_0) \ , \tag{1}$$

where $Pr(x_1, \ldots, x_N \mid x_0)$ is the probability of population compliance pattern $x_1, \ldots, x_N$ given policy $x_0$. The problem of inferring population responses from the available data is therefore a major part of the decision-making process.

## 2.2   Inference from empirical results for cohorts

Data cannot usually be collected about the response of every user to every possible policy choice. It is therefore clear that two major steps in the formation of a model-based decision process for security policy should be: the collection of relevant data on cohorts of users, and the process of generalisation to results about the entire population by inference.

Ideally, one would wish to infer (by Bayesian or statistical methods [47]) information about probability distributions over the spread of user behaviour. For example, one would wish to infer the distribution of $Pr(N.f(x_0) = m \mid x_0)$ for all integers $m$ and all policies $x_0$, or more generally $Pr(x_1, \ldots, x_N \mid x_0)$, as in the example from above. This may be hard, and assuming the distribution is reasonably peaked, one may settle for the most likely spread. The real problem here is in the collection of sufficient data.

A number of studies of security-related user behaviour have appeared in the academic literature. Zviran and Haga [71] studied just under one thousand users, finding information about password length and structure, frequency of change, frequency of write-down, memorability and correlation with sensitivity and importance. Yan, Anderson, Blackewll and Grant [69, 70] made a study of just under three hundred participants to produce results relating password memorability and security). Adams, Sasse and Lunt have carried out studies regarding usability of IT systems with given authentication mechanisms [2, 1]. Grawemeyer and Johnson [31] made a diary study of a small cohort of (twenty-two) users, that gave results about correlation between perception of required security and perception of required password strength, and between perception of perceived password strength and password structure. These existing studies appear to be efforts to gather data about user behaviour that will be broadly applicable to any organisation. This is of course very important. However, what is of interest in this paper is the kind of data that can be (quickly) gathered to support a given policy decision.

Real user behaviour is much richer than a simple binary secure-insecure response. Users typically have a multiplicity of many-part responses to policies. Worse still, the many parts are not usually chosen from a simple schedule, and users may not respond honestly to direct questions relating to non-compliance. The problem is exacerbated by the often multi-dimensional nature of the data required, variablitiy of behaviour across role, group, geography and function, and the fact that the significant risk may be posed by small numbers of extreme behaviours that are hard to detect.

Consider, for a moment, a detailed model in which the behaviour of each user matters to the firm. For simplicity, suppose that the experimenter constrains each user's response to $b$ binary responses, and that users respond honestly. That is, each user makes one of $2^b$ responses. If there are $M$ members of the cohort then there are $(2^b)^M = 2^{bM}$ possible responses by the cohort, but only $M$ actual responses gathered. It is likely (even in the case $b = 1$) that if the frequencies for the cohort are extrapolated to the population then most probabilities $Pr(x_1, \ldots, x_N \mid x_0)$ will be small, and very many will be estimated as zero. It may be that there are such zero probability estimates for which the corresponding loss $\mathcal{L}_0(x_0, x_1, \ldots, x_N)$ is very large. There could then be

large errors in the estimate of Equation (1). That is, the mix of a large, structured population, and heavy losses in rare circumstances make the expected losses hard to estimate. Now suppose that the population structure does not matter in the model, but only the overall frequencies of certain behaviours. A loss function of the form $\mathcal{L}_0(x_0, m)$ might then appropriate, and the problem is to infer $Pr(N.f(x_0) = m)$ for each integer $m$ and policy choice $x_0$. It may well be the case that this can be done for the existing policy choice $x_0$, but users may have difficulty knowing how they would respond to hypothetical policies. It is not clear that sufficient reliable data can be gathered in a timely way.

It is certainly infeasible to collect complete data for every policy decision on what changes in user behaviour would result: even where such effects are recorded in event logs they are only relevant to precise circumstances and only for one possible history (containing those policy choices actually made, attacks actually suffered and costs actually borne). Counterfactuals are invisible. One could perhaps study the effects of different policy choices on user behaviour in the laboratory, but it is likely hard to reproduce conditions facing users in their workplace and the effects of combination with existing policy, job requirements and social structure.

Regular, large scale data gathering exercises from user interviews and observations appear infeasible for two reasons. First, the constraints on time and on the number of available researchers. Second, organisations are unwilling to allow large-scale studies. Again, there seem to be two possible reasons: many organisations do not currently care about security enough for them to regard such exercises as justifiable in terms of cost (including disruption); alternatively some organisations are so sensitive to security that the introduction of researchers into their system itself constitutes an unjustifiable risk. A lack of access may be reflective of an implicit decision by the organisation about the value of security, but may be an artifact of the decentralised decision-making and the division of responsibilty for risk from that of other activities, say, IT operations (and revenue-generating functions).

The problem of understanding and guiding a user community is a problem of the social sciences. The difficulty of finding reliable and precise invariants of individual and population behaviour is a major source of uncertainty. It is therefore not surprising that it should be hard to accurately predict user responses to policy changes, except in very special cases. As a consequence it is extremely difficult to formulate a decision-making process for policy that is both practical and near-optimal. One will usually have to adopt heuristic methods.

## 2.3   Simulation

Computer simulation is a widely used technique for modelling complex systems [62, 22, 21]. In science it is of most use where convincing analytical models of a situation do not exist. However, it can also have a role as a source of intuitions and hypotheses and as a way of communicating model data (both inputs and results). The use of simulation may yield very precise answers (as in classical Monte Carlo methods [53, 52, 36, 47]) or be more heuristic. Their use in economics was pioneered by Simon [61, 62] as a way of getting to grips with the complex behaviour of bounded-rational agents. Further useful discussions of behavioural models in economics are found in [4, 41].

Simulation models can be constructed to represent security trade-offs and the behaviour of users [6, 8, 5, 9, 16, 10, 60]. This has several advantages. Models can be constructed relatively easily, and can be easily altered as assumptions are changed. Models can be built at varying levels of abstraction: detail is easy to include compared to equational formulations. The structure of models and the way that data is incorporated are relatively easy to communicate to stakeholders in decisions. All of these points are essential for practical support of decision-making. However, behind the second lies a further point. In systems with complicated behaviour, it is very difficult (and sometimes impossible) to know which details are relevant to the trajectory followed. It is therefore important to be able to study systems from many points-of-view.

The validity of results produced by simulation are often questioned because of the complexity of models, the difficulty of comprehending their execution path, and the apparent ease with which models can be adjusted to produce any desired result. The question of validation (or even falsification) for simulation models is a hard one, in part because they are applied where empirical results
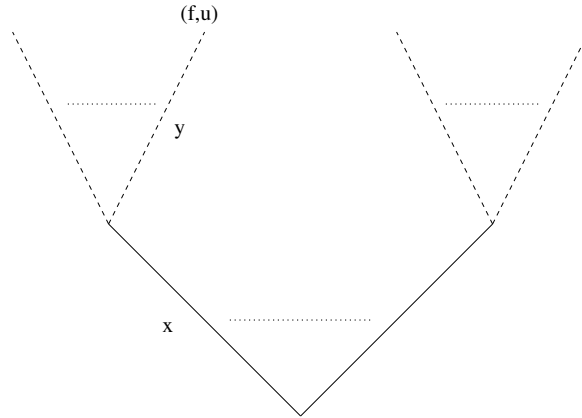
Figure 1: Policy setting as a Stackelberg game

and analytical results are not available, and sometimes because the natural phenomena under study are themselves fundamentally complex and difficult to test. Perhaps the only guidelines to protect against 'bad' models are discipline in the construction phase, and benchmarking in tests where results are available from another source.

For examples of the form outlined in Subsection 2.1 it is perfectly possible to construct simulation models in which both the population reaction $f(x)$, and the risk $r(x)$ arising from policy $x$ emerge from runs of the model. That is, the forms of these functions are generated by more primitive descriptions of processes, and the functions themselves are not required to be known in advance.

Many of the security simulation models (for example [6]) use hypothetical data to populate models: that is, they postulate plausible relationships or 'transfer curves'. These are grounded in the available data and on expert intuition. Of course, this is far from a perfect substitute for real data. The results of such simulations then must be understood as being contingent upon the accuracy of these assumption and the structure of the model. The models can be revised as both data and understanding of the scenario improve. The combination of small empirical studies, transfer curves and simulation models give heuristics for the decision-maker that have been demonstrated to be useful in practice.

## 3   The individual user as a rational agent

### 3.1   A Basic Framework

In the previous section users were described as though community behaviour was a simple function of the policy choice of the firm. However, user behaviour can also be derived from their preferences.

As a baseline model (for this section) the interaction of any firm and an individual user is a *Stackelberg game* (or follow-the-leader game) [11, 25]. In this game, the firm moves first by making a choice of policy, and the user responds with a reaction to that policy. In extensive form this can be viewed as a tree. A sketch of such a tree is shown in Figure 1. The policy choices have been drawn as solid lines and the user reactions as dotted lines (the fact that not all action choices have been shown is indicated by the horizontal dots). In this figure the firm receives payoff $f$ and the user $u$ if the firm chooses policy $x$ and the user responds with reaction $y$.

The user is assumed to be rational in the sense that he always chooses the highest available payoff, given the policy constraint set by the firm. Similarly, the firm wishes to produce the highest payoff for itself. Hence, assuming that the firm knows the form of the game tree and all of the payoffs, it can use backward induction to find optimal policies. That is, it can eliminate branches corresponding to user actions which would never be chosen becuase they are dominated

by branches that lead to greater payoff for the user. The calculation for the firm is then reduced to a simple utility calculation over its available policies.

## 3.2  Threshold behaviour

A simple illustration of how the Stackelberg framework could shed light on individual behaviour can now be given. This illustration is for a purely hypothetical firm and user — it is intended to be simple rather than realistic. There is no data to support the particular payoffs involved and the security problem is rather oversimplified. Nevertheless, the example shows how a security-related behaviour could be accounted for starting from individual preferences.

It is claimed in [7] that users typically exhibit a kind of threshold behaviour in their response to policy. The focus is on situations where the firm can set security policies, but the users have the option to comply or not. The idea is that the user has a maximum amount of effort (the *compliance budget*) that he is willing to expend on security related tasks. It seems rather obvious that elements of such behaviour exist in real user populations and, indeed, the results presented in [7] suggest that this is a non-negligible effect.

Consider a world in which a user faces a single security task. Let the threshold for the user be $t \geq 0$, the benefits of complying be $b \geq 0$ and the costs of complying be $c \geq 0$, all in the same units. The budget idea is encapsulated in a behavioural rule (adapted from [7]):

$$\text{comply if } c - b < t \text{ and don't otherwise}$$

This rule assumes that expected punishment from non-compliance is negligible (or counted as a benefit from compliance). Furthermore, in [7] attention is drawn to the approach to the threshold $t$ as ever more security tasks are required of the user — the 'hassle factor'. Thus it is really the marginal benefits and costs associated with the additional task (over and above those already accumulated) that need to be considered in the behavioural rule.

For simplicity, consider the case where no previous tasks have been assigned. The decision situation below is changed a little so that users have multinomial decisions rather than simple binary compliance. This allows richer threshold structures to emerge.

Imagine a stylised situation where a user is an employee of a firm, and is required to use a password in order to go about its work. The firm sets a policy for the length $n \geq 1$ of the password. The user responds by choosing a password of length $0 \leq m \leq n$. Let us ignore the fact that length of password is easy for the firm to enforce technologically, so that we can study the trade-offs implicit in the policy.

The choices for the payoffs chosen below are purely illustrative, but these choices are not unmotivated. For this exercise, it seems appropriate that the notion of cost to the individual user involves a term that reflects the mental load associated with remembering the password. Assume that this grows slowly initially with $m$, but then more-and-more rapidly: for further simplicity let it be quadratic. Issues to do with entropy and the ability of the human mind to make associations for longer passphrases are ignored. Now suppose that the firm imposes a linear punishment regime on the user, but that no further costs or benefits are involved in the decision. The user thus has a payoff of the form

$$\Pi_U(n, m) = -am^2 - b(n - m)$$

where the firm has chosen $n$ and the user $m$ for some constants $a$ and $b$. The fact that constants are used here reflects a simplification: user preferences are being treated as deterministic (perhaps through the use of expectation values) rather than as stochastic processes.

Simple calculus shows us that the user's optimal reaction is approximately $m = b/2a$ where $n \geq b/2a$, and $m = n$ otherwise. For $a = 1$ and $b = 19$ this line is shown in red in Figure 2. The actual optimal choices (which must be integers) are shown as blue dots. From Figure 2 one can see that a hard threshold has emerged in user behaviour.

The payoff $\Pi_0$ for the firm precisely is not given, as no optimization for the firm is to be carried-out here (although it would certainly be possible). However, assume that $\Pi_0$ is a strictly decreasing function of $m$, reflecting risk resulting from compromise of the password (which is assumed to be
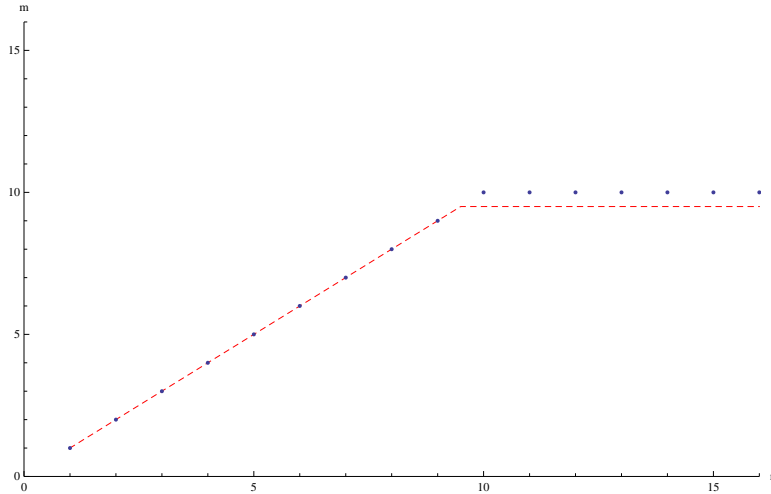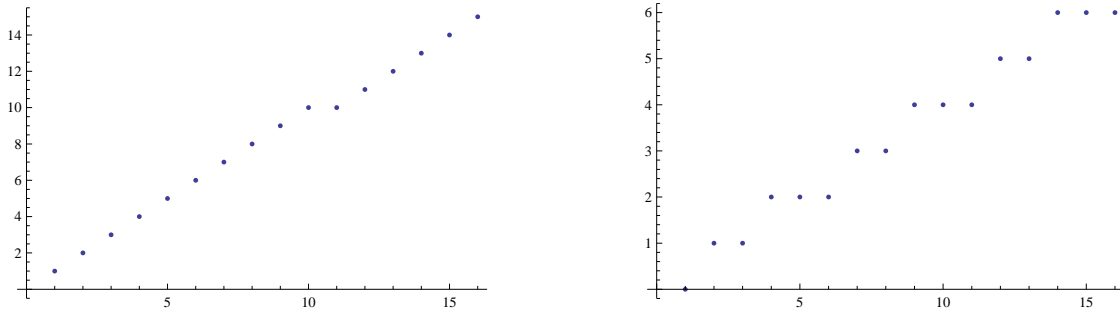
Figure 2: Threshold under linear punishment



Figure 3: Thresholds under quadratic punishment

less-and-less likely for increasing $m$). Suppose also that the function $\Pi_0$ is a decreasing function of $n$, say, because support costs associated with longer passwords increase. It is never worth the firm setting a policy $n > \text{ceiling}(b/2a)$ since the user can never be driven beyond this point.

Now consider the case where the user is subject to a quadratic punsihment regime, that is

$$\Pi_U(n, m) = -am^2 - b(n - m)^2$$

for all $m$, $n$. Calculus shows the optimal reaction for the user to a policy $n$ is approximately $bn/(a+b)$. Using the values $a = 1$, $b = 19$ gives the graph of user reaction on the left of Figure 3. In this graph there is a small stagger in the line of reactions given by the blue dots. As the ratio $a/b$ becomes larger so do the staggers, as shown in the right-hand graph for the values $a = 3$, $b = 2$.

The user now has a staggered threshold for compliance. With the same assumptions on $\Pi_0$ as before, it is never worth the firm's while investing in $n$ that leads to a user reaction on a stagger that lies on a plateau but not at the left-hand end. Under both punishment regimes the threshold behaviour has emerged from the combined preferences of the firm and the user.

## 3.3 Discussion of the framework

The Stackelberg framework given above is rather simplistic. The picture surrounding real decision-making is much more complex. Changes to policy will often be driven by changing user behaviour.

9

It is by no means clear that the firm moves first. The picture is further muddied by the presence of a changing threat environment featuring active agents and possible social-engineering attacks. As in Section 2 the subtleties introduced by time have been ignored. It has been assumed that the firm has complete information about the user — it knows his preferences exactly (and so can apply backward induction). It has also been assumed that the user has complete information about the firm (it is able to say know how hard it will be punished by the firm) and perfect information (it has internalised all policy). Neither of these assumptions is particularly realistic.

Perhaps the main stumbling block for the use of game theory as a predictive modelling tool is the difficulty of producing the loss function for the user without first observing behaviour. For very specific situations, it may be appropriate to do artificial experiments in the lab to discover user preferences, but this is surely not feasible in general. Discussion of a predictive variant of game theory that is rooted in Bayesian inference is taken up in Subsection 5.5.

A further very serious objection to the above type of model is the assumption that the behaviour of users can be modelled independently. Interaction is surely the key to much of their behaviour — people are after all social beings. For example, consider a user who is also an insider with malicious intent. His attack may only be possible because the firm assumes his preferences are those of a 'typical' user, when in fact they are chosen to be quite different. Thus his behaviour depends upon the firm's anticipation of other users. The interaction of users is the topic of Sections 4 and  5.

More technically sophisticated — and more detailed — games for explaining security phenomena have appeared in the literature. Nevertheless, the Stackelberg game remains a useful notion for framing arguments about the responses of typical users.

# 4   Multiple users with social interaction

## 4.1   From Micro to Macro

Security policy in large organisations can be complex. For a user to familiarise themselves with all its detail may take a considerable amount of time and effort. It is far from usual for a person to make decisions that are purely rational evaluations of trade-offs, say, between risk and productivity. More typically, they also assimilate with the prevailing attitudes that they perceive in others — in other words, they are strongly influenced by *social norms*.

From the point-of-view of policy decision-making, a firm will usually reason about the aggregate behaviour of the user community, rather than that of each individual, for the reasons explained in Section 2. Theories of how cultural macro-effects arise from individual micro-preferences *and* organisational and social structure would be extremely helpful.

In this section, the focus is upon a community of users that care about just two factors: first, maximizing revenue, and second, conforming (or not) with their colleagues. Different colleagues exert different amounts of influence over each user. The social effects that arise from such preferences can be extremely complex, so the focus here will be on simple binary decisions. Suppose that the firm's choice is either to adopt or not a policy in which a new security control is mandatory. Each of the users has a binary choice of compliance. Situations in which the users and the firm have multiple or a continuum of actions available can also be handled — the real source of the complexity in this model is in the interaction between colleagues.

The models considered will be dynamic: users are free to choose to act securely (comply) at some times, but act insecurely (not comply) at others. For simplicity the time horizon is unbounded. Furthermore, there is no cost to the user in switching their decision in either direction, at any time. In addition, the focus is upon the users' trade-offs between costs arising from compliance and the desire to conform (or not). The models produced relate to situations where it is difficult, undesirable or infeasible for the firm to attempt to find and act against any non-complying user, either by restricting access or some other punishment.

## 4.2 Preference and high-level interaction

The terminology of game theory is helpful for the description of the models, although they are not quite formalised as games here. There are $N + 1$ players: the firm and $N$ users. Let each $i$th player make a choice $x_i \in \{-1, 1\}$. The value $x_0 = 1$ means that the firm chooses to adopt the new policy, whilst $x_0 = -1$ means that it does not. For the $i$th user $x_i = 1$ means that he chooses to comply with policy, and $x_i = -1$ means that he does not. A *profile* is a tuple $\mathbf{x} = (x_0, \ldots, x_N)$. Such a profile may be written as $\mathbf{x} = (x_i, \mathbf{x}_{(i)})$ to emphasise the choice of the $i$th player: the $N$-tuple $\mathbf{x}_{(i)}$ is formed by removing the $i$th coordinate. An $N$-tuple of the form $(x_1, \ldots, x_N)$ is referred to as a *configuration*.

For the $i$th user and each $\mathbf{x}$ associate a loss of the form

$$\mathcal{L}_i(x_i, x_{(i)}) = -[h_i(x_i) + w_i(\mathbf{x})] \tag{2}$$

for some real-valued functions $h_i$ and $w_i$, where the index $j$ runs between 1 and $n$. Given $x_0$, the user wishes to minimize this loss, which is also known as the *local field on $i$ at $\mathbf{x}$*.

The function $h_i$ represents the preference of the user in the absence of social effects — it is called the *external field on $i$*. For a user $i$ who experiences a tension between the security policy and its preference, it will be the case that $h_i(1) < h_i(-1)$: if there were no social effects on $i$ ($w_i$ is identically zero) then the user would prefer not to apply the control. The function $w_i$ represents the influence of colleagues over the $i$th user — call it the *internal field on $i$*. The case where $h_i$ is identically zero (and the local field is the internal field on $i$) is known as the *free field on $i$*.

As a further simplfication the external and internal fields for each user $i$ are defined by

$$h_i(\mathbf{x}) = h_i x_i \qquad\qquad w_i(\mathbf{x}) = \sum_{j \neq i} w_{ij} x_i x_j$$

for constants $h_i > 0$ and $w_{ij} \in \mathbb{R}$ for $1 \leq j \leq n$. The constant

$$b_i = h_i + \sum_{j \neq i} w_{ij} x_j \tag{3}$$

is the local field for $i$ given $\mathbf{x}_{(i)}$.

The constant $w_{ij}$ represents the weight that the $i$th user attaches to agreeing (disagreeing) with the $j$th user. If $w_{ij} = 0$ then the $i$th user does not care what the $j$th user is doing. If $w_{ij} > 0$ then the term $w_{ij} x_i x_j$ influences the $i$th user towards agreement with the $j$th user. If $w_{ij} < 0$ then the term $w_{ij} x_i x_j$ tends to make the $i$th user disagree with the $j$th user. It is supposed that $w_{ii} = 0$ for all $i$. Note that it is possible that $w_{ij} \neq w_{ji}$, since it is usually not true that the mutual influence exerted by users is symmetric in strength.

Fix real constants $\beta_1, \ldots, \beta_n$ (these will be explained later). For each $\mathbf{x}$ define the *potential*

$$E(\mathbf{x}) \quad = \quad \sum_i \beta_i \mathcal{L}_i(\mathbf{x}) \tag{4}$$

for each tuple $\mathbf{x}$, with the index $i$ running between 1 and $N$. This is a kind of social welfare function, but note that it does not usually represent the preferences of the firm.

The interaction works as follows. First the firm acts (by choosing the policy), then the users follow. This is a little like the Stackelberg game of Section 3, but the game has a dynamic sub-game consisting of interacting users resulting from each policy choice. The firm acts just once, but the users act many times, adjusting their state if they find that they could have a lower loss. Note that this is not a single action — they may flip many times. The discussion of exactly how users interact is specific to each version of the model: examples are given in Subsection 4.4.

The choices user choices are really a function of time. Write $x_i(t)$ for the choice of user $i$ at time $t \geq 0$, and $\mathbf{x}(t) = (x_0, x_1(t), \ldots, x_n(t))$. The paramater $t$ will be dropped when the context allows. The details of how $x_i(0)$ is chosen will vary from model to model: for policies which arrive when a new technology is made available to users it will be appropriate to assume that all $x_i(0) = -1$ for all $i$; for policies which cover the use of an existing technology or business practice

this may not be the case. Let $P(t) = (x_1(t), \ldots, x_n(t))$ be the configuration at time $t$, and $P$ be the corresponding path through the space of configurations.

The preference of the firm should encapsulate both the cost associated with adopting the policy and the cost (risk) of users not conforming with the policy.

Some choices of the interaction between users are entirely detrministic, and the path $P$ is uniquely determined by the initial state $\mathbf{x}(0)$. In this case the loss for the firm takes the form

$$\mathcal{L}_0(x_0, x_1(0), \ldots, x_t(0)) = -[h_0(x_0) + h^0(P)]$$

for some functions $h_0$ and $h^0$ with $h_0(-1) < h_0(1)$ and $h^0(P) \geq 0$ for all $P$.

Suppose that for any $t$ there is an instantaneous loss function for the firm $\mathcal{L}_0^t(x_0, x_1(t), \ldots, x_n(t))$. In continuous time, the quantity $h^0(P)$ could be written as an integral [3] :

$$\int e^{-\gamma t} \ \mathcal{L}_0^t(x_0, x_1(t), \ldots, x_n(t)) \ dt$$

given a discount factor $\gamma$.

More generally, it is useful to consider models where the path arises as the realization of a stochastic process determined by $\mathbf{x}(0)$ and a collection of probability distributions. Assuming a probability measure $\mathbf{P}$ over the set of paths thus determined, the expected value of the loss to the firm takes the form of an expected loss over paths:

$$\mathcal{L}_0(x_0, x_1(0), \ldots, x_t(0)) = - \left[ h_0(x_0) + \int e^{-\gamma t} \ \mathcal{L}_0^t(x_0, P(t)) \ d\mathbf{P} \right]$$

using a suitable integral.

No optimization is to be performed here, so the form of the loss function is not examined in any more detail. The problem of understanding the dynamics is very hard, and this — and its impact on general loss functions of this form — is the focus of attention.

## 4.3 Spin models

The models of social effects that are going to be explored are closely related to models from statistical physics. *Spin models* [19, 20, 47] are used to explain how coupling between physical micro-components can lead to very important macro-properties of materials. Many of these, including the *Ising model*, are used to explain how large geometrical arrays of small magnets (particles, with binary values called spin), align (or do not) with neighbours with the same spin. Such models can then be shown to exhibit similar properties at the macro-scale as ferromagnetic materials. Perhaps the most impostant of these properties is the existence of a *phase transition*. Below a critical temperature the material is magnetic, because at the micro-level the spins agree (within regions called domains). Above the critical temperature, the spins become randomly arranged, and at the macro-level the magnetism $M$ diappears (in the absence of an external field). This is illustrated in Figure 4.

The description of social interaction given in Subsections 4.1 and 4.2 essentially constitute a spin model, although of a slightly unconventional kind. Probably the earliest social spin model is [24]. The same kind of model also arises in neural networks and learning theory. In Hopfield networks the firing of each individual neuron is related to the firing of certain neighbours. The potential (4) represents the propensity of the neurons to fire. Suitable choices of weights $w_{ij}$ allows for networks to be trained to reach chosen equilibria and thereby achieve various tasks.

The physical model of the dynamics of a spin model is driven by the potential (4) of the system. In the case of physical models all $\beta_i$ are equal. Let $\beta_i = 1$ for all $i$, and define the *inverse temperature* $\beta = 1/k_B T$, where $k_B$ is the Boltzmann constant and $T$ is the temperature. The *Gibbs measure* (sometimes called the Boltzmann distribution) is derived from the principles of

---

[3]This could be replaced by a discounted sum for a discrete time model.

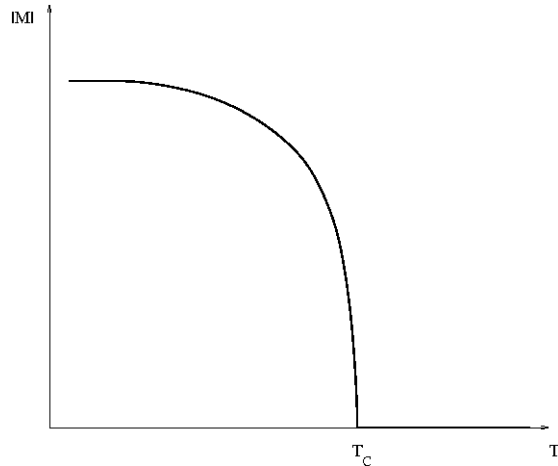Figure 4: Spontaneous magnetisation

statistical mechanics and gives the probability that in thermal equilibrium the configuration is $\mathbf{x}$ at inverse temperature $T$, given all the constants $h_i$ and $w_{ij}$. The measure is given by:

$$Pr(\mathbf{x}) = \frac{1}{Z} \exp\left(-\beta E(\mathbf{x})\right) \tag{5}$$

where $E(\mathbf{x})$ is the potential from (4), with all $\beta_i = 1$, and $Z$ the normalization constant. The constant $Z$ is called the *partition function* when it is treated as a function of the system paramaters.

Most macro-properties of spin systems are determined by the partition function. In physical systems, one is often interested in properties that occur as $N$ becomes very large. However, the partition function can be hard to calculate explicitly, except in special cases. Note that with $N$ spins there are $2^N$ states, so direct calculation of $Z$ becomes computationally infeasible for moderately large $N$. However, Monte Carlo simulation can be used to provide extremely precise estimates of the probability (5).

Near a phase transition many macro-properties (like heat capactiy) of a spin model are determined by a small number of *relevant parameters* (for example, the dimension of the geometrical lattice and the temperature) and not by *irrelevant paramaters* (for example, the exact details of the weights that couple spins). This is referred to as *universality*. It means that, even where specific interaction strenghts cannot be known, overall macro-properties can be accurately determined from the relevant parameters.

## 4.4 Dynamics of social models

User behaviour in reality is likeley determined by two intertwined processes, observation-and-reaction and anticipation-and-inference. Dynamic models are the natual way to model behaviour under observation-and-reaction. The play of games by tâtonnement and the checkerboard model of Schelling [58] are examples of this approach. More is said on anticipation in Section 5. For now, suppose that user behaviour is dominated by observation-and-reaction.

In order to produce a purely dynamical and completely accurate model, one needs a great deal of information about how (and when) users (jointly) update their state. It would be highly unusual for such information to be available to a decision-maker.

For simplicity, assume that users update asynchronously, and based only on the current configuration. Since real users may well communicate their intended behaviour with one another, purely asynchronous updating based on the current configuration may not be realistic. For example, if some policy changes, then it is not unlikely that users will communicate and choose approximately

13

simultaneously. Below we consider just two regimes, deterministic and Gibbs updating, but many others are possible.

### 4.4.1 Deterministic updating

Consider a regime in which users make a deterministic choice of action (comply or don't) after examining the state of their colleagues. Suppose that time advances in discrete steps. Assume that updating works by the repeated uniform random choice of a total ordering on the set of users. Thus, after updating, no user may update again until all other users have done an update. That is, updating takes place in rounds. If the update causes a change to the sign of a spin then this is referred to as a *flip*.

Consider, for example, the following toy example of physical security. Suppose that the employees of a firm go to the canteen every day, and in the same groupings of colleagues. Suppose that the canteen is not a secure area, but that the employees work in a secure area. On the way back, every employee must pass through a door which is only wide enough to allow one person to pass at once. The firm may worry that it should institue a pass system on the door, so that the door opens and closes automatically, but only immediately after a swipe from a security pass that each employee carries. Suppose that under the proposed system it is possible for employees to tailgate, that is, for multiple employees to enter sequentially from a single swipe before the door closes. If it went ahead with the system the firm would have a policy that users should not tailgate. The firm does not want to bear the cost of additional guards to ensure that the policy is not violated or violators punished. Suppose that each employee makes his choice by observing the most recent choices of every member of his group. In reality, the first employee does not have the opportunity to tailgate, but let us suppose that he can somehow pass through the door without authentication.

Each individual wishes to minimize their loss funtion $\mathcal{L}_i$ from Equation (2). Therefore if it is the $i$th users' turn to update, then he assumes that the other users behaviour $\mathbf{x}_{(i)}$ is fixed, and he tries to choose $x_i$ to minimize $\mathcal{L}_i(x_i, \mathbf{x}_{(i)})$. All he has to do is compare the losses at the $x_i = 1$ and $x_i = -1$. Note that $\mathcal{L}_i(x_i = 1, x_{(i)}) - \mathcal{L}_i(x_i = -1, x_{(i)}) = -2(h_i + \sum_j w_{ij} x_j)$.

Let $\Theta$ be the step-function:

$$\Theta(y) = \left\{ \begin{array}{ll} 1 & \text{if } y \geq 0 \\ -1 & \text{if } y < 0 \end{array} \right. . \tag{6}$$

The update $x_i(t+1)$ to the $i$th user (spin) $x_i(t)$ is:

$$x_i(t+1) = \Theta(h_i + \sum_j w_{ij} x_j) \tag{7}$$

where $x_j = x_j(t)$ if the $j$th spin has not yet been updated in this round, and $x_j = x_j(t+1)$ if it has.

In the special case where $h_i = 0$ for all $i$, and where the $w_{ij}$ are symmetric this makes the system an *asynchronous, binary Hopfield network* [47] The presence of asymmetric weights greatly complicates a problem. Consider a system with just two users. User 1 desires to be like User 2, with weight $w_{12} = 1$, but User 2 likes to be different to User 1 with weight $w_{21} = -1$, and $h_1 = h_2 = 0$. Flips then occur indefinitely far into the future [4] .

Deterministic models of human communities are of limited use, because the behaviour of members in complicated choice situations is typically best described as stochastic. Nevertheless, the fact that Hopfield nets can be *constructed* to achieve certain tasks (like pattern recognition) by choice of couplings suggests that communities could similarly be formed to achieve tasks, if only the couplings could be controlled.

---

[4]Almost surely, in the probabilistic sense [33])

#### 4.4.2 Random utility and logistic response

In laboratory experiments human subjects have been shown to produce stochastic, rather than deterministic, responses to choice situations. Theories of preference that accurately account for real behaviour must therefore often have a stochastic component. The corresponding utility theory is known as *random utility theory* [12, 44, 45, 50, 48].

Suppose that each user's loss function (2) is perturbed by a random term $\epsilon_i$ that is a function of the users' choice. That is, re-define $\mathcal{L}_i(\mathbf{x})$ by:

$$\mathcal{L}_i(x_i, \mathbf{x}_{(i)}) = -[h_i(\mathbf{x}) + w_i(\mathbf{x})] - \epsilon_i(x_i) = -[(h_i + \sum_j w_{ij}x_j)x_i + \epsilon_i(x_i)] = -[b_i x_i + \epsilon_i(x_i)]. \quad (8)$$

Given $\mathbf{x}_{(i)}$, the user prefers $x_i = 1$ if and only if $\mathcal{L}_i(1, \mathbf{x}_{(i)}) < \mathcal{L}_i(-1, \mathbf{x}_{(i)})$, and so if and only if $\epsilon_i(-1) - \epsilon_i(1) > 2b_i$.

A common assumption for such a binary choice is that the random variable $\epsilon_i(-1) - \epsilon_i(1)$ has the *logistic distribution* with mean 0 and scale paramater $2/\beta_i$:

$$\Pr\left(\epsilon_i(-1) - \epsilon_i(1) \le z\right) = \frac{1}{1 + \exp\left(-\frac{1}{2}\beta_i z\right)}. \quad (9)$$

The $i$th user then chooses $x_i$ with probability:

$$\Pr\left(x_i \mid \mathbf{x}_{(i)}\right) = \frac{1}{1 + \exp\left(-\beta_i b_i x_i\right)}. \quad (10)$$

Call the paramaters $\beta_1, \ldots, \beta_N$ *local inverse temperatures*. They are sometimes viewed as capturing the level of rationality of the agents. The update formula (10) encapsulates the fact that $\beta_i$ captures the $i$th user's level of certainty about how to react. For fixed $\mathbf{x}_{(i)}$, as $\beta_i \to 0$, $Pr(x_i = 1 \mid \mathbf{x}_{(i)}) \to 1/2$, so the user is very uncertain about how best to respond. As $\beta_i \to \infty$, $Pr(x_i = 1 \mid \mathbf{x}_{(i)}) \to 1$ if $b_i > 0$, and $Pr(x_i = 1 \mid \mathbf{x}_{(i)}) \to 0$ if $b_i < 0$.

Thus user repsonses to their colleagues current states can be made stochastic. This might happen for two reasons: the user makes a rational calculation (of $b_i$) but then deliberately randomizes his response; the user's calculation of $b_i$ contains errors (because of some form of bounded rationality).

#### 4.4.3 Gibbs updating

Recall the situation described for deterministic updating of user behaviour. Now suppose that the user updating process is uncertain, either because the users themselves are stochastic entities, or because their true updating process is opaque but can be approximated by a distribution.

Suppose that $\mathbf{x} = \mathbf{x}(t)$ is given. Users update in rounds as with deterministic updating. Choose some user $i$. The quantity $b_i$ is calculated using the most recent value of $\mathbf{x}_{(i)}$ (using those components already updated in the round). The new update rule for each $i$th user is that given by Equation (10). That is:

set $x_i(t+1) = 1$ with probability $(1 + \exp\left(-\beta_i b_i\right))^{-1}$, and $x_i(t+1) = -1$ otherwise.

Remove the integer $i$ from the set $1 \le i \le N$, choose the next variable to update, and continue until the round is complete. This is referred to as *Gibbs updating*.

This updating process is the same thing as *Gibbs sampling* [47] upon the probability distribution (5) with $E(x)$ given by (4) and $\beta = 1$, and under the assumption that the users' randomization processes (the samples of Equation (10) for different $i$) are independent. Note that statistical correlations between users will still arise because of the coefficients $w_{ij}$. The updating process constructs Markov chains $(\mathbf{x}(t) \mid t \in \mathbb{N})$, which can be combined with Monte Carlo methods to approximate the equilibrium distribution (5).

In physical models such approximation is important because it is usually not feasible to calculate the partition function $Z$ directly, but Gibbs sampling eventually converges arbitrarily close to the distribution (5) [5].

Turning this around, given our model of user behaviour, if run for sufficiently long, it will, on average, converge to the distribution (5), with (4) and $\beta = 1$. To be explicit, the probability that the system is in configuration $\mathbf{x}$ is:

$$Pr(\mathbf{x}) = \frac{1}{Z} \exp\left(-\sum_i \beta_i \mathcal{L}_i(\mathbf{x})\right) .$$

Note that, even though the users are not playing as a team and no agent is trying to optimize $E$ as social welfare, this potential is still at the heart of the dynamics of the system.

## 4.5 Social macro-effects from micro-behaviour

The significance of spin models in physics is in accounting for the way in which macro-effects come about through the dynamics of coupled micro-behaviour, including phase changes in ferromagnetic materials. Phase changes are radical changes to properties of physical systems that occur as the result of small changes to system paramaters (often temperature). A familiar example is the phase transition of water to ice (and vice-versa) at freezing point. Magnetic materials exhibit a phase transition at a critical temperature called the Curie point (or temperature). Below the Curie point ferromagnets can be magnetised; above it they cannot. At the micro-level of spins, this is accounted for by spins becoming uncorrelated because of additional energy at higher temperature.

In the Ising spin model one has a macro-quantity called the *magnetization*, at any instant given from the state $\mathbf{x}$ by,

$$M = \mu \sum_{i=1}^{N} x_i ,$$

where $\mu$ is a real constant that gives the magnetization of a single spin. The phase change is illustrated in the sketch of average (modulus) magnetization against temperature shown in Figure 4.

Define an instantaneous loss function at time $t$ for the firm by a weighted combination of the user behaviours:

$$\mathcal{L}_0^t(\mathbf{x}) = -\sum_{i=1}^{N} \mu_i x_i$$

for $\mathbf{x} = (x_0, x_1, \ldots, x_N)$ and some constants $\mu_1, \ldots, \mu_N$, and $\mathbf{x}$ is the state at time $t$.

The instantaneous loss of the firm could reasonably considered to be an economic generalization of the notion of magnetization. Bearing in mind the phase change in magnetization around changes in a system paramater, it seems essential to know whether or not a phase change (or other non-linearity) occurs in the instantaneous loss function of the firm.

In the Ising model the phase change occurs at a critical value of the temperature paramater, or equivalently the *inverse temperature* $\beta$. In the Gibbs updating model the question is whether there is a phase change around critical constellations of the tuple of local inverse temperatures $\beta_1, \ldots, \beta_N$. For the sake of simplicty it may sometimes be useful to consider the same change of scale to each of the local inverse temperatures. That is, each $\beta_i$ may be written in the form $\beta \beta^i$ for some $\beta^i$, and only the paramater $\beta$ is varied.

Typically, a firm can influence the local inverse temperatures: it can raise or lower the uncertainty of its users about the level of advantage in compliance. For example, it could randomly perturb the costs of compliance (the $h_i$ components) according to some scheme for which the users cannot easily guess the distribution. However, note that local inverse temperatures are not typically under the complete control of the firm. They are influenced by perceptions about the stability of the firm, or, for a security example, of how the threat environment potentially affects

---

[5]Except in certain pathological cases.

users losses. Note that changes to local inverse temperatures made by the firm have costs that are not accounted for in the present loss functions. However, the question about phase changes still stands.

Without detailed empirical study it is impossible to know whether phase changes really exist in social populations, and whether spin models can properly account for their behaviour. Nevertheless, the following examples may be worthy of study as plausible candidates.

1. File-sharing of copyrighted material within a single jurisdiction. Suppose a standardised punishment. An individual could choose to participate or not. Broader society plays the role of the firm. There is coupling between individuals through social interaction both on- and offline.

   Society could decide to punish sharers and reward non-sharers in an random way, for which it is difficult for the user to estimate the distribution (either in scale or likelihood of successfull prosecution). This roughly corresponds to adding uncertainty to the $h_i$ components. Alternatively, society could introduce unpredicatable noise into the reputation and feedback mechanisms in the forums used by sharers. This would make the coupling constants $w_{ij}$ uncertain. If a phase change change occurs in sharing behaviour with respect to the level of disruption then society has the potential to capitalize upon this highly nonlinear response. As with magnetization, as the $\beta$ paramater approaches the cititical value the behaviour of individuals suddenly becomes uncorrelated. It may be cheaper for society to achieve this kind of limited result than it is to try to wipe-out sharing altogether through a constant change to the punishment regime (that is, the external field components $h_i$).

   Similar phenomena could occur in other online communal activities that are considered antisocial by broader society: criminal activities (including credit-card and bank-detail fraud, unlawful pornography), terrorist activities, and certain cults. Note also the possibility for an attack by one social community on another (or the provider of the service infrastructure for that community on another) by similar means.

2. Now consider a population that can choose to use (or not) some legitimate online service (e.g. banking). An unpredicatably evolving threat environment could be represented as setting a high $\beta$ paramater (the external field components $h_i$ fluctuate unpredictably). In particular, it would then be difficult for individuals in the population to estimate the distribution underlying social and technological shocks from threats. If there is a phase change then, as the critical value of $\beta$ is approached the demand for the service would rapidly reduce to almost zero, as would the derived social welfare.

3. Consider a situation in which the majority of employees of a firm do not conform with a policy of the firm. This might be: a 'no-tailgaiting' policy on sites; mandatory encryption of all company data stored on portable media; a strict policy that the firm's IT resources are only used for its business.

   The addition of unpredicatable punishments and rewards (by the firm) to the external field components $h_i$ would correspond to raising the uncertainty $\beta$. For sufficiently high $\beta$ the employees' choices would become almost uniformly random. Again, this may be cheaper than changing the external field in a fixed way.

More generally, suppose that a firm wishes to control a community that is a reasonably modelled as a spin system. Suppose that the community is initially in an undesirable state (lots of non-compliance). The firm could change rewards and punishments predictably (change the external field) to try to drive the community into a more desirable equilibrium. This has an analogue in physical spin systems — it is the attempt to change the sign of the magnetization $M$. The picture of the (free) energy associated with such spin systems is shown in Figure 5. The activation energy $E^*$ required to force the change of sign in $M$ is typically large. Similarly, the costs for the organisation may be prohibitively large. Howver, if the spin system has a phase transition for a critical uncertainty $\beta_c$ then a three-stage alternative exists that may be cheaper. First, alter
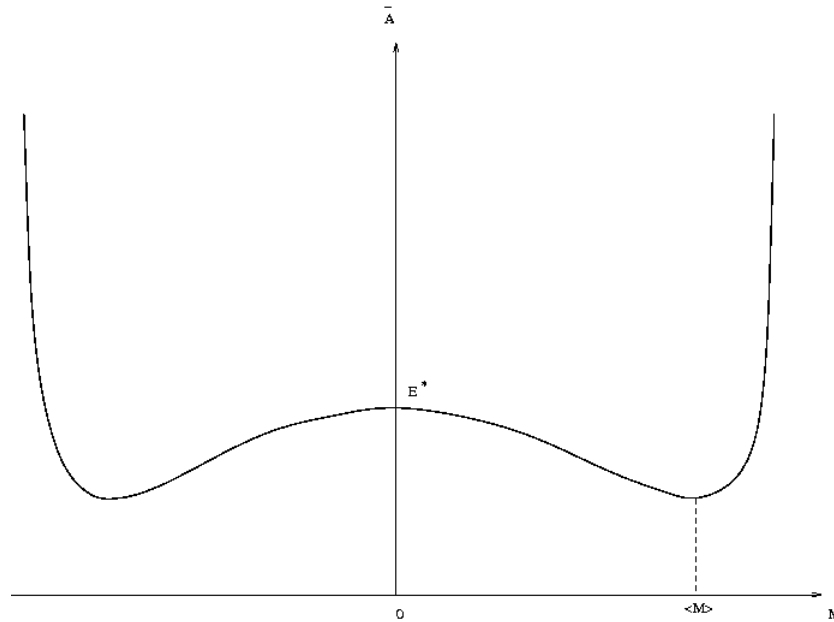
Figure 5: Reversible work function in a free field

the uncertainty of each employees about which is their best state, thus changing the $\beta$ paramater until it is close to $\beta_c$. This could be done either with unpredictable rewards and punishments, or by disrupting the couplings between community members. The correlation between community members is greatly reduced at this point. Second, apply a predictable punishment (and/or reward) system (this is the external field) to drive the system to the desired social loss (corresponding to magnetization). Third, as optimal, reduce the punishment/reward system. The language (unfreezing-transition-crystallization) of Lewin's theory of change management in organisations [42, 27, 46] is extremely suggestive of this type of process.

The third stage may be possible due to the phenomenon of *symmetry-breaking* in spin systems — that is, their tendency not to return to a state of zero average spin value once the external field is removed. At the macro-level, in the case of magnetization, the familiar concept of *hysteresis* is observed, where a ferromagnetic material (below the Curie temperature) may start in a state of zero-magnetization, become magnetized in the presence of an applied external field, but retain a (reduced) magntization after the field is removed. This is shown in Figure 6. Thus hysteresis and symmetry breaking could be exploited to appropriately set the organisation's long-term punishment and reward strategy. Each of the above stages would have costs, so of course, deciding how exactly to carry-out the above three-stage process would depend on each of these.

## 4.6 Comparison of physical and social spin models

There are a number of important differences between physical spin models and the social models proposed above.

The potential (4) can be regarded as an instantaneous social loss function. Note, however, that this neither usually corresponds directly to the firm's loss, nor does it account immediately for the behaviour of users, since each one only cares about minimizing its own loss. Therefore there is no immediate justification to use the potential to drive the dynamics of the social model. The use of the potential can be justified when the users update according to a Markov Chain Monte Carlo sampling of (5). More general updating strategies, which are surely more natural in most cases, need not lead to this analysis. Blume has shown [13] that various forms of updating using
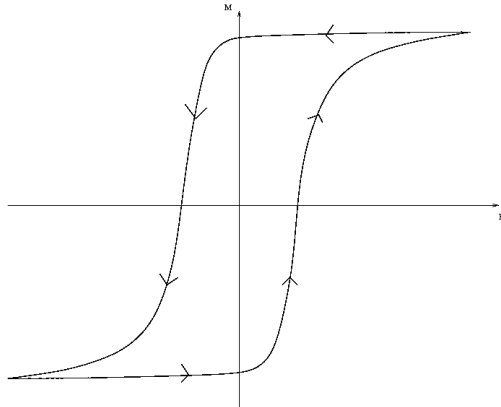
Figure 6: Hysteresis

players with 'Poisson alarm clocks' have stationary distributions which have support that can be characterised by standard game-theoretic equilibria. In general it is far from obvious that any analytic form for the probability measure over states exists. If this is the case, or the partition function is difficult to evaluate, then simulation and statistical analysis will be the only viable approach.

It is not the case that every social spin system subject to an arbitrary updating regime reaches a static equilibrium. The appropriate notion of equilibrium is that of stationary distribution. In Monte Carlo simulations of physical spin systems one typically only cares about the long run behaviour, after the transient effects of initial conditions have been eliminated. However, in the security applications suggested here, in which a firm sets a policy, and the community follows different considerations may need to be made. The initial conditions may well be known, at least to a good approximation, for example when a firm changes policy. The firm will, in general, care about effects of the policy change that occur before, as well as after, the community approximately stabilizes (assuming it does). Such transient effects cannot be eliminated unless they are small. Thus, even when the equilibrium distribution eventually approximates the Gibbs measure, it may be that other distributions are required for evaluating the expected losses before equilibrium. There may be serious consequences for an organization that ignores major disruption caused by a change in security policy, even if the user community eventually settles into a desirable equilibrium.

The most important use of spin models in physical systems is in accounting for phase changes. It is not obvious that any given spin system will have phase changes — on the other hand there may be many, because there are many different coupling constants (weakly coupled sub-networks in the user community would have phase-changes at lower temperatures than strongly-coupled networks). It is known that the existence of phase changes can depend upon the topolgy of coupling, for example, the one-dimensional Ising model does not have a phase change, but the two-dimensional model does. The underlying topology required for studying social effects is very different from that in geometrical spin models. For example, in a large corporation, it may well be appropriate to suppose that the individuals in some department are primarily informed by the behaviour of their colleagues and direct managers. The structure of the social and organizational network (as encoded by the weight parameters $w_{ij}$) determine social effects in a critical way. In addition, it is usually not the case that the weights are all uniform or even symmetric. It is fashionable to model social networks as small-world graphs which are characterised by the fact that most vertices are not directly connected, but are connected by a small number of edges [65, 66, 55]. Spin models based on such graphs have been studied [37]. In [26] instances of phase changes in communities of artificial agents playing a team game as a spin system are given.

The exact value of $\beta$ at which phase changes occur is, in general, difficult to predict in physical spin models. However, it can also be found by simulation given the exact valuse of all the constants.

In physical models the way that certain macro-parameters change about any phase change often can be predicted from a small number of relevant paramaters (for example, by the dimension of the lattice of spins), whereas much of the detail of individual couplings is unnecessary. This is an instance of universality. It would be of interest to know what kind of universality properties there are in social and organizational networks, since for example, the exact vales of all the individuals coupling constants cannot be known. See also [63] for a discussion of potential instances of universality in networks of artificial agents.

In general, the range of behaviour possible by users will be much more complex than a simple binary choice. If users make a choice from multiple quantal values then the Potts generalization of the Ising model is an evident model to start from. Above, pairwise multiplications have been used to compare spins in the users' loss functions, but the pairwise comparisons could be quite different. User preferences may be even more complicated, requiring the consideration of not just bipartite comparisons, but tripartite and beyond. Such situations require consideration of further generalised individual user losses.

In real organisations the user population changes over time. One might suggest an analogy with techniques like those of grand canonical ensembles in statistical physics for models of such situations. However, there is a further generalization that one should really consider at the same time, namely the fact that the structure (as well as the size) of organizations changes over time. In a spin-model, this corresponds to coupling constants strengthening and weakening over time, possibly reversing their sign, and even fluctuating stochastically. Again, simulation will often be the only useful modelling methodology.

In the physical spin systems it is often appropriate to assume that all of the coupling strenghts have the same known, fixed value, and that the influence of an external field on all spins is only dependent on the magnitude of the field and the sign of the spin. In a model of a real social system with a large number of users, it is clearly utterly infeasible to establish all of the coupling strengths and individual effects of punishment and reward. It is also infeasible to establish all the uncertainty paramaters $\beta_i$. Thus, whilst spin models have explanatory powers over cultural effects in user communities, for example why one sub-group complies with a policy but another does not, it is unclear to what extent such models can be used to predict the effects of changes made by the organization, and what the costs might be. Together with the question above about universality, the most important unanswered question here is which approximations are most appropriate, and for modelling which type of social situations. When can users be treated as though they were decoupled? When can they be treated as though all couplings were of the same strength? When can couplings strengths be drawn from some distribution? Which updating strategies are useful approximations to the actual fluctuations that the community can make?

## 5    Mutliple coupled users and simultaneous decision

### 5.1    The choice situation

In the previous section the models of user communities were dynamic and users' choices to comply with policy were able to fluctuate backward-and-forward over time. This section will be about irreversible binary compliance decisions that are taken simultaneously by users (at least, without conferring or mutual observation). Real situations in which every user's behaviour is purely anticipative must be rather rare. Clearly though, one needs to have a grasp on such situations before one can comprehend more complicated mixed models of both anticipative and reactive behaviour. The following, slightly artificial, examples are of the type treated in this section.

1. A population of employees of a firm separately attending a virtual meeting. Employees may or may not conform with some set of policies to ensure that their connection is secure. Employees have beliefs about how their colleagues behave gathered from previous social contact.

2. A firm becomes aware that a particular application has a security flaw, but a patch is

expected shortly. It issues a policy that the application should not be used until the patch is applied. The firm also has a strict policy that users should not discuss with colleagues any aspect of their behaviour that relates to security. Users may or may not choose to use the application, and not to do so may harm their productivity.

In each situation for the users, there is a cost in complying, a benefit in agreeing with their colleagues, but no opportunity to observe those colleagues

This section focusses on the behaviour of the user community after the policy has been set. As a matter of notational convenience the ambient policy choice of the firm ($x_0$ in the previous sections) is omitted from the loss fucntion.

## 5.2   User behaviour from belief about colleague behaviour

Brock and Durlauf have advocated the use of statistical mechanics as a methodology for dealing with interactions between agents in economic models [23, 14]. This is referred to below as the *BD-method*. Their approach is based on pure anticipation of colleague behaviour, rather than the dynamics of observed behaviour. This subsection contains a model in their style, although the details are changed in an insignificant way to make it more similar to the dynamic models above.

There are $N$ users each of whom as a choice to comply [6] . This is represented by $x_i \in \{-1, 1\}$ for each $1 \le i \le N$. It is additionally supposed that each user $i$ has a given subjective expectation value $\mathbf{E}^i(x_j)$ for the behaviour $x_j$ of each other user $j \ne i$.

As a consequence each local field $b_i$ (as in (3)) has an expected value:

$$\mathbf{E}^i(b_i) = h_i + \sum_{j \ne i} w_{ij} \, \mathbf{E}^i(x_j) \ .$$

Agent behaviour is determined by minimisation of a loss function $L_i(x_i)$ over the two alternatives. This is defined by

$$L_i(x_i) = \mathcal{L}_i(x_i, \mathbf{E}^i(\mathbf{x}_{(i)})) \ ,$$

where $\mathcal{L}_i$ is the loss function from (8), $\mathbf{E}^i(\mathbf{x}_{(i)})$ is the $N - 1$ tuple with $j$th component $\mathbf{E}^i(x_j)$ for $j < i$ and $\mathbf{E}^i(x_{j+1})$ for $j \ge i$. Therefore $L_i(x_i) = -[\mathbf{E}^i(b_i)x_i + \epsilon_i(x_i)]$ for some random variable $\epsilon_i(x_i)$. Note that the loss function is assumed to be random. Suppose that the random variables $\epsilon_i(-1) - \epsilon(1)$ have the distribution given in (9). This allows for the derivation of the conditional probability

$$Pr(x_i \mid h_i, \mathbf{E}^i(x_j) \forall j \ne i) \propto \exp\left(-\beta_i \, \mathbf{E}^i(b_i)x_i\right) \tag{11}$$

for each agent $i$, since $x_i$ is not chosen if $L_i(x_i) > L_i(-x_i)$.

Assuming that agents choose simultaneously and independently, but based on their subjective beliefs, the conditional probability of the joint choice $\mathbf{x}$ is

$$Pr(\mathbf{x} \mid h_i, \mathbf{E}^i(x_j) \forall i, \forall j \ne i) = \frac{1}{Z} \prod_i \exp\left(-\beta_i \, \mathbf{E}^i(b_i)x_i\right) = \frac{1}{Z} \exp\left(-\sum_i \beta_i \, \mathbf{E}^i(b_i)x_i\right)$$

where $\mathbf{x}$ is the vector of agent choices, and $Z$ is the normalization constant.

Thus the Brock-Durlauf method leads to a Gibbs measure for the distribution over states $\mathbf{x}$, using $\beta = 1$ and the potential

$$E(\mathbf{x}) = \sum_i \beta_i \, \mathbf{E}^i(b_i)x_i \ .$$

If each $\epsilon_i(x_i)$ has a subjective expectation of zero, that is $\mathbf{E}^i(\epsilon_i(x_i)) = 0$), then

$$E(\mathbf{x}) = \sum_i \beta_i L_i(x_i) = \sum_i \beta_i \mathcal{L}_i(x_i, \mathbf{E}^i(\mathbf{x}_{(i)})) \ .$$

Notice that the dynamics have been simplified through the use of (subjective) expectation of colleague action. The fixing of expected colleague behaviour (spin) fixes the local field (expected loss for each choice) for each user. The more general case of fixing the local field without fixing user belief about colleagues is considered in Subsection 5.3.

---

[6]The word comply is used in the same sense as in the earlier sections.

## 5.3   Mean field theory: behaviour under deviation from typical loss

Suppose that each individual user $i$ has a belief $\mathbf{E}^i(b_i)$ about his local field. The user does not necessarily have a belief about the behaviour of the other users. To this end, suppose that the loss function can be written as

$$\mathcal{L}_i(\mathbf{x}) = -[\mathbf{E}^i(b_i x_i) + \epsilon_i(x_i)]$$

for all $\mathbf{x}$, with $b_i$ as in (3) and $\epsilon_i$ satisfying the same conditions as in (9). Note that the BD-method is a special case. Then:

$$Pr(x_i \mid \mathbf{E}^i(b_i)) \propto \exp\left(-\beta_i\, \mathbf{E}^i(b_i)x_i\right)$$

for some normalization constant $Z_i$. Therefore $i$ has

$$\mathbf{E}^i(x_i) = \tanh\left(\beta_i\, \mathbf{E}^i(b_i)\right)$$

as the expectation for the state of its own spin

Suppose that each user's evaluation of the local field is accurate. That is $\mathbf{E}^i(b_i) = \mathbf{E}(b_i)$ for all $i$. This is rather like a *rational expectations* assumption. Then $i$'s subjective expectation for the value of its spin agrees with the mathematical expectation: $\mathbf{E}^i(x_i) = \mathbf{E}(x_i)$ . The set-up above is essentially equivalent to the *mean field approximation* of a physical system. From the point-of-view of every spin, fluctuations of the other spins have no impact upon it as individuals, but rather only through their average.

Let

$$\Delta b_i = \sum_{j \neq i} w_{ij}\, \mathbf{E}(x_j)$$

so that $\mathbf{E}(b_i) = h_i + \Delta b_i$. The expected state of the $i$th user is then

$$\mathbf{E}(x_i) = \tanh\left(\beta_i(h_i + \Delta b_i)\right) . \tag{12}$$

This gives a system of $N$ simultaneous equations with unknowns $\mathbf{E}(x_1), \ldots, \mathbf{E}(x_N)$. This may have solutions for certain constellations of paramaters $\beta_1, \ldots, \beta_N$.

In the physical case, one is usually only interested in the situation with a single type of spin, so that all the constants $\beta_i$ are equal, and so that all the expectation terms in $\Delta b_i$ are equal. One finds that the mean magnetization per spin $m = \mathbf{E}(x_1)$ must be given by solutions to the equation

$$m = \tanh\left(\beta h + \beta z w m\right) \tag{13}$$

where $\beta = \beta_i$, $h = h_i$, $w$ is the coupling constant for neighbouring spins, and $z$ is the (uniform) number of neighbours for any spin. In the free field ($h_i = 0$) the single resulting equation has non-zero solutions for all values of $\beta$ below, and no values above, a critical value $\beta_c$, since $\beta$ determines the gradient of the right-hand-side of (13) at the origin. This is illustrated in Figure 7. The existence of a non-zero solution corresponds to the possibility of spontaneous magnetization, and $\beta_c$ is the inverse temperature of the Curie point $T_c$ at which the phase change occurs. For a low value of $\beta$ one has the lower curve of Figure 7, with only a trivial solution to the single equation.

In order to apply this result directly to an economic scenario, one would have to assume that all agents were identical, and were arranged so that they all had the same number of neighbours. The conclusions about the structure of equilbria found in [14] follow this pattern. It does not seem easy to generalize the argument to account for heterogeneous agents in more realistic networks.

Note also that mean field theory is an approximation. It produces the misleading prediction of a phase transtion for the $1-D$ Ising model, and numerically inaccurate results for the critical temperature $T_C$ in the $2-D$ case (although it fares better in higher-dimensions). It is often the case that economic situations played by simple reaction are shown to converge to the same equilibria as would be found by mutual anticipation [11, 25]. The failures of the mean field approximation show that there are networks of agents and values of $\beta$ for which eventual reactive behaviour will not converge to the anticipative equilibria of the BD- or mean-field methods.
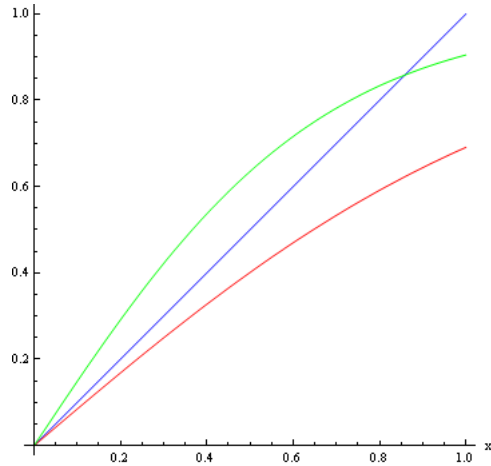
Figure 7: Mean field solutions

There are other ways of proving the existence of phase changes in spin models that give more accurate results: for example variational mean field theories and renormalization group theory (for simple presentations see [47, 19] respectively). Doubtless these can be translated into economic language to give results about equilibria in communities of economic agents. These methods (and the various standard calculations of the partition function) use the fact that spins are arranged in a geometrical lattice. It is entirely unclear whether any of these methods can be modified to apply to more realistic networks of heterogeneous agents.

## 5.4   Quantal games: behaviour from belief about colleague preferences

In the BD-method models, the $i$th user chooses independently, conditioned on a belief that the behaviour of other agents has a fixed expected level. Each such fixed level $\mathbf{E}^i(x_j)$ does not depend upon the actual behaviour $x_j$ selected by the $j$th agent, nor upon its preferences. In game-theoretic terms the level of anticipation is very simple. Of course, there may also be anticipation by an agent trying to control the community by setting conditions in the form of paramaters for the BD-model. Instead, consider a situation in which each agent is aware that its colleagues choices may be conditional on their beliefs about it, and in which each agent has a fixed level of belief about its expected loss under every strategy combination, and in which all of these losses are public knowledge (for the community). Evidently, there is the possibility for each agent to exploit this additional knowledge, that is to engage in properly game-theoretic behaviour. The notion of *quantal response* games (here called QR-games) [51] generalize the BD-method to such situations.

For example, consider a situation where there is a user community who work at home most of the time. The community is organised into workgroups. Each user's preferences are determined solely from those of his workgroup. Each user is resposible for the maintenance of a client laptop. The workgroup meets only periodically. An instruction to apply some emergency patch or not to engage in some newly-declared insecure activity is issued by email to the users. Making this change would cause disruption to each user, but each user also desires to be like his colleagues. Users cannot observe each other and do not confer because of time-pressure and because their communications are monitored. Each user must try to anticipate his colleagues' behaviours using his knowledge of their preferences, but knows that they will also try to anticipate him.

In the standard analysis of an $N$-player (strategic form) game, it is usual to derive from the payoffs a *reaction function* that describes how each player would best respond to each profile of strategy choices by all the other players: an equilibrium is then a strategy profile which is a suitable fixed-point for the product of reactions.

QR-games give a more flexible framework for describing real behaviour of human players. They

are based on the idea that players cannot exactly calculate the payoff that they would receive (given a strategy profile for the other players), and instead the payoff is known only up to some (continuous) probability distribution. Thus any player responds (to a given strategy profile for the other players) as though they are sampling a probability distribution (over their pure strategies). That is, some noise is introduced into the players' calculations, and this results in uncertainty in their play. For the $i$th player, with available pure strategies in the set $X_i$, there is a *statistical reaction function* $r_i$ which arises from the distribution over payoffs. For any profile $\mathbf{x}_{(i)}$ for the other players, the value $r_i(\mathbf{x}_{(i)})$ is a probability distribution over $X_i$: write $r_i(\mathbf{x}_{(i)})(x_i)$ for the probability that the strategy $x_i$ is chosen given $\mathbf{x}_{(i)}$, that is $Pr(x_i \mid x_{(i)}) = r_i(\mathbf{x}_{(i)})(x_i)$. Under mild conditions on the game and on the payoff distribution, it is guaranteed that there is at least one equilibrium in such reactions, this is known as a *quantal response equilibrium* (QRE).

Let $\mathcal{L}_i(x_i, \mathbf{x}_{(i)})$ be the loss for the $i$th player under the strategy profile $(x_i, \mathbf{x}_{(i)})$. It is assumed (in particular) that

$$\mathcal{L}_i(\mathbf{x}) = -b_i(\mathbf{x}) + \epsilon_i(x_i) \tag{14}$$

where $b_i$ is a deterministic function, and the random variable $\epsilon_i(x_i)$ is distributed so that $\mathbf{E}(\epsilon_i(x_i)) = 0$. Note that the loss of (8) is a special case. For simplicity, assume rational expectation for every use, so that

$$\mathbf{E}^i(\mathcal{L}_i(\mathbf{x})) = \mathbf{E}(\mathcal{L}_i(\mathbf{x})) = -b_i(\mathbf{x}) \ .$$

The expectation of $i$ of its losses under any strategy profile $\mathbf{x}$ is known.

The logistic response to any $\mathbf{x}_{(i)}$ is

$$r_i(\mathbf{x}_{(i)})(x_i) = \frac{1}{Z_i} \exp(-\beta_i \, \mathbf{E}(\mathcal{L}(x_i, \mathbf{x}_{(i)}))) \ , \tag{15}$$

for some constant $\beta_i$ and normalization constant $Z_i = \sum_{x_i \in X_i} \exp(-\beta_i \, \mathbf{E}(\mathcal{L}_i(x_i, \mathbf{x}_{(i)})))$.

Players in simple strategic form games have been observed to have statistical reaction functions [51]. The parameters $\beta_i$ that determine the sharpness of the logistic statistical response, can be estimated from experiments by logistic regression. However, these paramaters can change as players learn to play the game.

In the case where each player has a binary choice and each $b_i$ takes the form given in (3) the logistic response is identical to the updating rule (10). Play of the QR-game by tâtonnement would thus be like Gibbs updating, but with simultaneous rather than asynchronous updates.

Any instance of a QRE in the logistic response case is known as a *Logit Equilibrium* (LE). Let $L_i(x_i)$ be the function with $L_i(x_i)(\mathbf{x}_{(i)}) = \mathcal{L}_i(x_i, \mathbf{x}_{(i)})$ for every profile $\mathbf{x}_{(i)}$. Each LE consists of a family of probability distributions $q = (q_1, \ldots, q_n)$, where each $q_i$ is a distribution over the strategy set $X_i$, satisfying

$$q_i(x_i) = \frac{1}{Z^i} \exp(-\beta_i \, \mathbf{E}_{q_{(i)}}(L_i(x_i))) \tag{16}$$

where $Z^i$ is the normalization constant, and $\mathbf{E}_{q_{(i)}}(-)$ is the expectation operator, taken with respect to the distribution $q_{(i)}$ across strategy profiles $\mathbf{x}_{(i)}$. Since the players sample these distributions indepndently the prediction across strategy profiles is

$$q(\mathbf{x}) \quad = \quad \prod_i q_i(x_i) = \tfrac{1}{Z} \exp\left(-E(\mathbf{x})\right) \tag{17}$$

with $E(\mathbf{x})$ as in (4), and $Z = \prod_i Z^i$.

Hence, the use of the Gibbs measure (with $\beta = 1$ and potential (4)) to predict the distribution over user behaviours is emerges from users making a single simultaneous move in a quantal game.

This result can, in principal, be applied to both small and large groups of users. The problem of collecting all model parameters for real players will become infeasible as the size of the population increases. It may be that the strategy sets of users can be taken to be equal (for example, for a binary compliance decision) and that the parameters $\beta_i$ can be chosen by drawing from an esitmated distribution across the population.
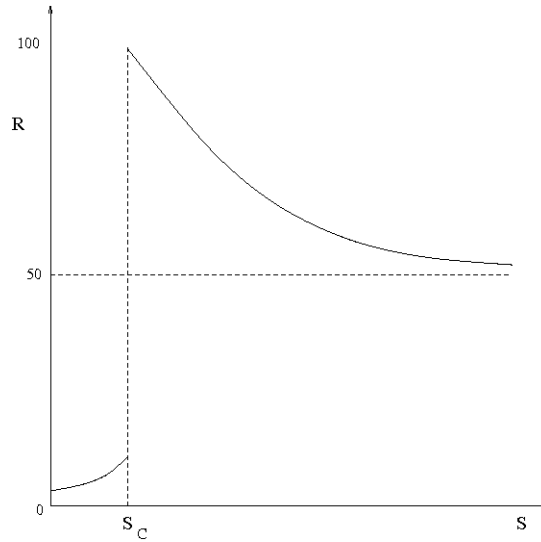
Figure 8: Threshold models: equilibrium frequency of rioters

## 5.5  Alternative approaches

Influential dynamic models of populations were produced by Schelling [58] in his studies of segeregation, and later extended to other phenomena. Granovetter [30] further extended this work to produce models of the behaviour of populations in which members have a threshold for participation in some activity (referred to, generically, as rioting). Each member chooses to riot as soon as he observes that his threshold number of members are already rioting. The interaction is dynamic. The models can be regarded as simple instances of Hopfield nets. Some security problems could be explained in these terms. The willingness of individuals to partcipate in online criminal markets may be a good example, since they may seek safety in numbers, regardless of social ties. For most organisations, and most communities of users, it will not be publically available information what the actual numbers of non-conformists with some policy are. Rather, each user may treat his colleagues as a sample, so simple numerical threshold models can be rather too simple.

Rioting behaviour is evidently critically dependent upon the structure of thresholds in the population. Granovetter considered what happens when this structure is created by sampling probability distributions. For a particular example with population members drawn from a certain normal distribution the results exhibit a sharp non-linearity. Namely, at a certain critical value of the standard deviation, participation jumps from around six percent of the population, to around one hundred percent. This is illustrated in Figure 8, with $R$ the percentage of rioters, $S$ the standard deviation, and $S_C$ the critical value. This is rather like a phase change. However, because of the structure of the model, it is quite different. First, note that all members are deterministic, so the population will not feature members flipping backwards and forwards between peaceful and rioting. Second, since the models are all deterministic, the uncertainty in behaviour comes from the sampling. The 'phase change' here is a flip from one highly-ordered state to another, rather than a flip between order and randomness as in a spin model.

Granovetter noted that populations in which certain members exhibit more influence over another are also of interest — this just means allowing for more general weights in the net. The practical difficulties for population members of accurately assessing whether their thresholds have been reached are mentioned — thus the threshold calculations could be peturbed by a random variable, as in the random loss models above. It is also noted that certain activities have maximum, as well as minimum thresholds. For example, one could imagine that a user will not deviate from a security policy if he does not see enough of his colleagues deviating, but will suffer if too many of his colleagues deviate, since for example, the over-arching organisation will face an attack or

bring in draconian measures to enforce policy. Such maximums are neither handled in threshold models nor in existing spin models. Examples of this type are discussed in work on the El Farol problem [4], including the results of simulations, and in the literature on minority games [18, 17]. They could perhaps be handled in economic spin models by augmenting the loss terms associated with each spin by adding to each user's loss terms representing the negative effects of too many colleagues agreeing to deviate.

Wolpert has initiated an entirely new view on game theory with an emphasis on the prediction of results of games. This work is summarised in [68, 67] and uses techniques from statistical physics and information theory. His starting point is the rigorous use of an extended Bayesian methodology — maximum entropy (maxent) methods, and the minimization of a loss function over the possible predictions. The prediction mechanism used by the modeller is constrained to this methodology. That is, the prediction (say, a joint mixed strategy $q$) should be the choice of the modeller which minimizes his given loss. The prediction is made by applying Bayesian updating to a given likelihood function (estimable by experiment) and a given prior. Thus it can be argued that the modeller is applying a properly scientific, rather than merely mathematical programme, and that his predictions are therefore better justified. Wolpert argues that the so-called entropic prior should be used, since one is trying to predict the probability $q$. This is in line with the principle that a prediction of a probability distribution should be chosen to maximize entropy (in the sense familiar to cryptographers), that is minimize information, subject to any constraints [39, 56, 59]. For a caveat on the application of maximum entropy see [47]. To be further explicit, in this view one does not assume that any prior notion of equilibrium should be the modeller's prediction. Rather, standard equilibrium notions can be shown to agree (or not) with the methodologically-approved predictions. The prediction that eventually appears is usually approximated by a Gibbs measure. This distribution is often similar to a QRE, specifically a LE, as in Subsection 5.4. However, in a QRE each players reaction does not take account of the other players' uncertainties. Wolpert's method does take this into account and produces corrections to the QRE. It also produces a resolution to the problem of prediction for games with multiple equilibria, and allows for variable numers of players. Although technically challenging, this methodology may be the most appropriate for predictive modelling of user groups with anticipatory behaviour, because of its design as an inference process to work on empirical data. However, even this methodology must surely struggle when faced by the need to gather data on large numbers of users.

Different probability models based on graphs [32] and interacting particle systems [43] are no doubt better suited to modelling other security situations involving user communities. For example, percolation models and contact processes are better suited to modelling the progress of the uptake of a new security technology in a structured community where, once a user adopts the technology, he will never go back (as with acceptance of most automated security updates).

# 6   Conclusions

Decisions by an organisation about security policy are often attempts to influence and control the behaviour of users of the system. Such decisions may have impact upon the well-being of the organisation and so should be supported by firm evidence. Ideally, the evidence should be generated by predictive models that are supported by carefully collected data. The appropriate type of model depends on the policy choices, the structural details of the organisation, its preferences, its context (including other rational agents), and the data available.

It is hard to come up with a practical methodology that meets this ideal. Many policy decisions are taken quickly in response to a rapidly changing environment. The entire decision-making process must be fast enough to provide timely results. There are significant difficulties in gathering sufficient data on the reaction of the user community to all possible policy choices, particularly to enable a direct calculation of expected loss (or utility). Such a data gathering exercise may be quite different to those carried-out without time constraints in a spirit of pure enquiry regarding general reactions of users to existing policies.

Part of the difficulty in predicting the responses of the community comes from the fact that

users are themselves rational agents, and may not perceive policy as being in their best interests. They may try to find strategies to work around policy or to balance its perceived cost against perceived risk. Such a situation — featuring autonomous agents following their own preferences — is naturally understood to be of the type treated by game theory. However, it seems that it would be very difficult to collect sufficient data on the preferences of users to make a practical methodology for decisions with time constraints. This is not to say that games do not provide useful, general insights into the behaviour of users.

A further difficulty in guaging the response of the user community lies in the interaction between users. Individual users are often strongly influenced by their colleagues. Interactions between large communities of agents can be highly complex. The complexity can arise from the rationality of the agents, and bounds on that rationality. However, complexity can also arise from interaction, and even relatively simple agents can give rise to unpredictable behaviour. The behaviour of a community is not usually described by simply aggregating the preferences of indivduals. Models of the behaviour of two-way interactions between the organisation and a typical user (or a distribution of users) are of limited use. One might contemplate the use of an $N+1$-player game (with $N$ users plus the organisation), but acquiring data and making predictions would surely be too difficult.

Spin models give a way to study macro-effects arising from interactions between (relatively simple) coupled agents. It seems reasonable that such models could describe user communities in simple security situations. It would be rather hard to gather data to completely populate such a model and to make accurate predictions. Rather, the value in such models is in alerting us to the possibility of sudden changes in community behaviour — phase changes. A phase change in a user community would (usually) result in a radical, non-linear change in the organisation's performance (measured by some loss function). A decision-maker would have to take such phenomena into account. The spin models suggest that phase changes could take place when users are subject to a critical level of uncertainty about their environment. Empirical studies on whether phase changes exist in real user populations would be much more significant than mathematical models. Producing such empirical evidence may not be easy; it could only really be done under controlled conditions. If phase changes are present then universality properties should be sought.

Detailed models (of organisation, users or loss) are difficult to use for prediction as they demand more data than can feasibly be gathered. Often, however, a significant level of detail is necessary in order not to produce misleading results. Rich models of systems can be hard to capture equationally. For a security decision-maker, simulations (possibly using hypothetical data) may be a valuable way to explore the potential consequences of organisational decisions about security. This is particularly true when the decision-maker faces strict time constraints. As always, the results of inferences from models need to be interpreted and evaluated with regard to the definitiveness of data, models and the inference process itself.

### Acknowledgements

# References

[1] A. Adams and M. A. Sasse. Users are not the enemy: Why users compromise security mechanisms and how to take remedial measures. *Communications of the ACM*, 42(12):40–46, 1999.

[2] A. Adams, M. A. Sasse, and P. Lunt. Making passwords secure and usable. In *Proceedings of HCI on People and Computers XII*, pages 1–19. Springer, 1997.

[3] R. Anderson. Why information security is hard: An economic perspective. In *Proc. 17th Annual Computer Security Conference*. IEEE Computer Society, 2001.

[4] W. B. Arthur. Bounded rationality and inductive behavior (the El Farol problem). *American Economic Review*, 84:406–411, 1994.

[5] A. Baldwin, M. Casassa Mont, and S. Shiu. Using modelling and simulation for policy decision support in identity management. In *Proceedings of the 2009 IEEE International Symposium on Policies for Distributed Systems and Networks*, pages 17–24. IEEE Computer Society, 2009.

[6] A. Beautement, R. Coles, J. Griffin, C. Ioannidis, B. Monahan, D. Pym, M. A. Sasse, and M. Wonham. Modelling the human and technological costs and benefits of USB memory stick security. In M. E. Johnson, editor, *Managing Information Risk and the Economics of Security*, pages 141–163. Springer, 2009.

[7] A. Beautement, M. A. Sasse, and M. Wonham. The compliance budget: managing security behaviour in organisations. In *Proc. 2008 Workshop on New Security Paradigms*, pages 47–58. ACM, 2009.

[8] Y. Beres, M. Casassa Mont, J. Griffin, and S. Shiu. Using security metrics coupled with predictive modeling and simulation to assess security processes. In *Proceedings of the 2009 3rd International Symposium on Empirical Software Engineering and Measurement*, pages 564–573. IEEE Computer Society, 2009.

[9] Y. Beres, J. Griffin, S. Shiu, M. Heitman, D. Markle, and P. Ventura. Analysing the performance of security solutions to reduce vulnerability exposure window. In *Proc. 2008 Annual Computer Security Applications Conference*, 2008.

[10] Y. Beres, D. Pym, and S. Shiu. Decision support for systems security investment. In *5th IFIP/IEEE International Workshop on Business-driven IT Management (BDIM 2010)*. IEEE Xplore Digital Library, 2010. To appear. Draft at: `http://www.cs.bath.ac.uk/~pym/BeresPymShiu.pdf`. Checked 12th Feb. 2010.

[11] K. Binmore. *Fun and Games: A Text on Game Theory*. D.C. Heath, 1992.

[12] H. D. Block and J. Marschak. Random orderings and stochastic theories of responses. In Olkin, Ghurye, Hoeffding, Madow, and Mann, editors, *Contributions to probability and statistics*. Stanford University Press, 1960.

[13] L. E. Blume. The statistical mechanics of strategic interaction. *Games and Economic Behaviour*, 5:387–424, 1993.

[14] W. Brock and S. Durlauf. Discrete choice with social interactions. *Review of Economic Studies*, 68(2):235–260, 2001.

[15] D. Burghes and A. Graham. *Introduction to Control Theory including Optimal Control*. Ellis Horwood Limited, second edition, 1986.

[16] M. Casassa Mont, Y. Beres, D. Pym, and S. Shiu. Economics of identity and access management: Providing decision support for investments. In *Proc. Business-driven Information Management (BDIM) 2010*. IEEE Xplore, 2010. To appear.

[17] D. Challet, M. Marsili, and G. Ottino. Shedding light on El Farol. *Physica A: Statistical Mechanics and its Applications*, 332:469–482, 2004.

[18] D. Challet and Y.-C. Zhang. Evolution of cooperation and organization in an evolutionary game. *Physica A: Statistical Mechanics and its Applications*, 246:407–418, 1997.

[19] D. Chandler. *Introduction to Modern Statistical Mechanics*. Oxford University Press, 1987.

[20] B. A. Cipra. An introduction to the Ising model. *The American Mathematical Monthly*, 94(10):937–959, 1987.

[21] M. Collinson, B. Monahan, and D. Pym. *A Discipline of Mathematical Systems Modelling*. College Publications, 2010. To appear.

[22] M. Collinson, B. Monahan, and D. Pym. Semantics for structured systems modelling and simulation. In *Proc. Simutools 2010 (to appear)*. ACM Digital Library; EU Digital Library, 2010.

[23] S. Durlauf. How can statistical mechanics contribute to social science. *Proceedings of the National Academy of Sciences of the United States of America*, 96(19):10582–10584, 1999.

[24] H. Föllmer. Random economies with many interacting agents. *J. Math. Econ.*, 2:51–62, 1974.

[25] D. Fudenberg and J. Tirole. *Game Theory*. MIT Press, 1991.

[26] R. Glinton, K. Sycara, D. Scerri, and P. Scerri. Statistical mechanics of belief sharing in multi agent systems. *Information Fusion Journal*, Special Issue on Agent Based Fusion, 2010. Available online.

[27] L. Goodstein and W. Warner Burke. Creating successful organization change. *Organizational Change*, 19(4):5–17, 1991. Reprinted in [46].

[28] L. A. Gordon and M. P. Loeb. The economics of information security investment. *ACM Transactions on Information and Systems Security*, 5(4):438–457, 2002.

[29] L. A. Gordon and M. P. Loeb. *Managing Cybersecurity Resources: A Cost-Benefit Analysis*. McGraw Hill, 2006.

[30] M. Granovetter. Threshold models of collective behaviour. *The American Journal of Sociology*, 83(6):1420–1443, 1978.

[31] B. Grawemeyer and H. Johnson. How secure is your password? Towards modelling human password creation. Proceedings of the The First Trust Economics Workshop, University College London, England 23 June 2009. Available at: `http://www.trust-economics.org/TEWorkshopProceedings.pdf`. Checked 3rd Feb. 2010.

[32] G. R. Grimmett. *Probability on Graphs*. Cambridge University Press, 2010. To appear. Draft at: `http://www.statslab.cam.ac.uk/~grg/books/USpgs.pdf`. Checked 21st Feb. 2010.

[33] G. R. Grimmett and D. R. Stirzaker. *Probability and Random Processes*. Oxford University Press, second edition, 1992.

[34] J. Grossklags, N. Christin, and J. Chuang. Security investment (failures) in five economic environments: A comparison of homogeneous and heterogeneous user agents. In *Proceedings (online) of the 7th Workshop on Economics of Information Security (WEIS 2008)*, 2008. Available at: `http://weis2008.econinfosec.org/papers/Grossklags.pdf`. Checked 19th Feb. 2010.

[35] L. Hassell and S. Wiedenbeck. Human factors and information security. Manuscript. Available at: `http://repository.binus.ac.id/content/A0334/A033461622.pdf`. Checked 21st Feb. 2010, 2004.

[36] W. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1):97–109, 1970.

[37] C. Hauert and G. Szabó. Game theory and physics. *The American Journal of Physics*, 73(5):405–414, 2004.

[38] C. Ioannidis, D. Pym, and J. Williams. Investments and trade-offs in the economics of information security. In *Proceedings of Financial Cryptography and Data Security '09*, volume 5628 of *LNCS*, pages 148–166. Springer, 2009.

[39] E. T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, 2003. Edited by G. L. Bretthorst.

[40] M. E. Johnson and X. Zhao. The value of escalation and incentives in managing information access. In M. E. Johnson, editor, *Managing Information Risk and the Economics of Security*. Springer, 2009.

[41] D. M. Kreps. *A Course in Microeconomic Theory*. Financial Times/ Prentice Hall, 1990.

[42] K. Lewin. Frontiers in group dynamics 1. *Human Relations*, 1:5–41, 1947.

[43] T. M. Liggett. *Interacting Particle Systems*. Springer, 1985.

[44] R. D. Luce. A probabilistic theory of utility. *Econometrica*, 26(2):193–224, 1958.

[45] R. D. Luce. *Individual Choice Behaviour*. Wiley, 1959.

[46] C. Mabey and B. Mayon-White. *Managing Change*. Paul Chapman Publishing Ltd., 2nd edition, 1993.

[47] D. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.

[48] T. Magnac. Logit models of individual choices. In *New Palgrave Dictionary of Economics*. Palgrave Macmillan, 2nd edition, 2009.

[49] R. E. Marks. Validating simulation models: a general framework and four applied examples. *Computational Economics*, 30(3):265–290, 2007.

[50] D. McFadden. Quantal choice analysis: a survey. *Annals of Economic and Social Measurement*, 5(4):363–390, 1976.

[51] R. McKelvey and T. Palfrey. Quantal response equilibria for normal form games. *Games and Economic Behaviour*, 10:6–38, 1995.

[52] N. Metropolis and E. Teller. Equation of state calculations by fast computing machines. *Journal of Chemical Physics*, 21(6):1087–1092, 1953.

[53] N. Metropolis and S. Ulam. The Monte Carlo method. *Journal of the American Statistical Association*, 44(247):335–341, 1949.

[54] D. A. Nadler. Concepts for the management of organizational change. In *Readings in the Management of Innovation*. Ballinger Publishing Company, 1980. Reprinted in [46].

[55] M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.

[56] J. B. Paris. *The Uncertain Reasoner's Companion: A Mathematical Perspective*. Cambridge University Press, 1994.

[57] E. M. Rogers. *Diffusion of Innovations*. Glencoe: Free Press, 1962.

[58] T. Schelling. Dynamic models of segregation. *Journal of Mathematical Sociology*, 1:143–186, 1971.

[59] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27:379–423,623–656, 1948.

[60] R. Shay and E. Bertino. A comprehensive simulation tool for the analysis of password policies. *International Journal of Information Security*, 8:275–289, 2009.

[61] H. A. Simon. A behavioural model of rational choice. *The Quarterly Journal of Economics*, 69(1):99–118, 1955.

[62] H. A. Simon. *The Sciences of the Artificial*. MIT Press, 3rd edition, 1996.

[63] H. Van Dyke Parunak, S. Brueckner, and R. Savit. Universality in multi-agent systems. In *Third International Joint Conference on Autonomous Agents and Multi-agent Systems (AAMAS'04)*, pages 930–937, 2004.

[64] H. Varian. System reliability and free riding. In *Economics of Information Security*, volume 12 of *Advances in Information Security*. Springer, 2004.

[65] D. Watts. *Small Worlds: The Dynamics of Networks between Order and Randomness*. Princeton University Press, 1999.

[66] D. Watts and S. Strogatz. Collective dynamics of 'small-world' networks. *(Letters to) Nature*, 393:440–442, 1998.

[67] D. Wolpert. A predictive theory of games. Manuscript, 2008. Available at: `http://arxiv.org/abs/nlin/0512015v1`. Checked 3rd Feb. 2010.

[68] D. Wolpert. Information theory — the bridge connecting bounded rational game theory and statistical physics. In D. Braha and Y. Bar-Yam, editors, *Complex Engineering Systems*. Perseus Books, 2004.

[69] J. Yan, A. Blackwell, R. Anderson, and A. Grant. The memorability and security of passwords — empirical results. Technical Report 500, University of Cambridge, 2000. Available at: `http://www.cl.cam.ac.uk/techreports/UCAM-CL-TR-500.pdf`. Checked 3rd Feb. 2010.

[70] J. Yan, A. Blackwell, R. Anderson, and A. Grant. Password memorability and security: empirical results. *IEEE Security and Privacy*, 2(5):25–31, 2004.

[71] M. Zviran and W. Haga. Password security: An empirical study. *Journal of Management Information Systems*, 15(4):161–186, 1999.