

Regulating Autonomous Agents facing Conflicting Objectives: A Command and Control Example

Jim Q. Smith and Lorraine Dodd
Warwick University and Cranfield University

Abstract

UK military commanders have a degree of devolved decision authority delegated from command and control (C2) regulators, and are trained and expected to act rationally and accountably. Recent experimental results suggest that experienced commanders usually appear to act as if they are subjective expected utility maximizers. The only scenarios where this appears not so is when the immediate mission objectives conflict with broader campaign objectives. Then the apparent rationality of even experienced commanders often evaporates. In this paper we show that if the C2 regulator assumes her commander is expected utility maximizing and that he uses a suitable multiattribute utility function then, even when she is remote from the field of action and her information is sparse, this regulator can nevertheless predict when scenarios might lead her commanders into making irrational decisions.

1 Introduction

To encourage its personnel to act as flexibly as possible, a command and control (C2) regulator (usually at a senior level of UK military command), devolves differing degrees of decision authority to field commanders. Within a mission statement, such a commander is given detailed information about how his decisive acts might provoke praise or inflame public opinion, promote or undermine trust and hence expedite or frustrate campaign objectives in the medium and long term. We henceforth say these factors frame the *acceptability* objective. Such factors tend to encourage restraint, appeasement, the avoidance of conflict and embrace the need for public security. This objective will be complemented by the commander's particular mission task, called his *immediate* objective,

which directs him to achieve certain immediate military aims against an adversary. The commander is trained and expected to act rationally. In particular, when given decision autonomy, he is aware that he will need to be able to justify his chosen acts in later after action reviews. In this paper, following the example of [7] a "rational" agent will be interpreted as someone who appears to act as a subjective expected utility (SEU) maximizer.

Surprisingly, in contrast with some other domains, the careful analysis of evidence from on-going studies of decision-making under simulated conditions of internal command contention and situational uncertainty, applied to the domain of military C2 on serving and experienced UK military commanders, see e.g. [1],[3],[8] have identified only one challenging scenario where the acts and explanations given by experienced commander were not consistent with those of an SEU decision maker. This was when there was no act that could score well on both the immediate and acceptability objectives and that simultaneously a compromise between two attributes scored poorly relative to more extreme alternatives. Then several commanders appeared to exhibit various irrational behaviors such as hypervigilance [6] or decision suppression [2]. Even when such agents managed to remain rational, such *bipolar* scenarios defined circumstance when discontinuity across adjacent autonomous commanders might be induced: see below. The object of this paper is therefore to study how and why such scenarios can arise and so inform the C2 regulator about how to train and communicate missions to its commanders to minimise the incidence of such scenarios. The challenge for the C2 regulator is that, even if she is prepared to accept the hypothesis that her commander acts like an SEU maximizer in non challenging circumstances, she is *remote* from the unfolding situation on the ground. She is therefore unlikely to know more than some broad qualitative features of her commander's utility function and posterior density. Despite this lack of precise information we demonstrate in this paper that the regulator can nevertheless still draw useful conclusions about how she should communicate missions and training which encourages behavior which is rational in the sense above.

Perhaps the simplest way for the C2 regulator to examine when bipolar scenarios might arise is to assume that her SEU maximizing commander has a utility with two attributes - one broadly scoring success in acceptability, the other scoring success in the immediate objectives associated with the mission. Such a model is described in the next section. Under this hypothesis to avoid bipolar situations C2 would then need to avoid facing her commander with situations where his expected utility might exhibit multiple isolated maxima. In Section 3 we illustrate how

a bipolar scenario can occur in a typical military operation. Reporting some new general results from singularity theory as they apply to this setting, [4], in Section 4 we then proceed to outline how, under certain mild regularity conditions, the remote regulator can use these models to address the effectiveness of the decision making of her agent more generally. The broad conclusion of the paper is that the C2 regulator should not ask much more from a commander than he thinks he could deliver in a sense made more precise below.

2 Some Technical Apparatus

A rational commander must discover the *intensity* d^* of action maximizing the expectation $\bar{U}(d)$ of his utility function $U(d, \mathbf{x})$. Here the regulator C2 assumes her commander's utility function $U(d, \mathbf{x})$ has two *value independent attributes* $\mathbf{x} = (x_1, x_2)$, [5], [11] where the attribute x_1 will reflect immediate factors and x_2 acceptability factors. Then

$$U(d, \mathbf{x}) = k_1 U_1(d, x_1) + k_2 U_2(d, x_2)$$

$d \in \mathbb{R}$. Recall that the *marginal utility* $U_i(d, x_i)$ - here a function of its argument only, taking values between 0 and 1 - scores the quality attainment of x_i using intensity d , $i = 1, 2$. The *criteria weights* k_i satisfy $k_i \geq 0$, $k_1 + k_2 = 1$. They reflect the relative importance of attaining high scores U_i , $i = 1, 2$ and allows the decision maker to trade off success between the two attributes appropriately.

The expected utility can then be written as

$$\bar{U}(d) = k_1 \bar{U}_1(d) + k_2 \bar{U}_2(d) \quad (1)$$

where, for $i = 1, 2$,

$$\bar{U}_i(d) = \int U_i(d, x_i) p_i(x_i) dx_i$$

and $p_i(x_i)$ denotes his posterior marginal density of x_i , $i = 1, 2$. Let $u_i^- = \inf_{d \in D} \bar{U}_i(d)$ and $u_i^+ = \sup_{d \in D} \bar{U}_i(d)$ denote, respectively, the worst and best possible expected scores under their respective posterior marginal densities on each of the attributes.

Note that because of the assumed linear form of the utility function, we do not need to assume the two attributes are stochastically independent of each other, provided that the decisions contingent on a choice of intensity associated with immediate and acceptability consequences respectively do not constrain one another. In [4] we argue that this second condition will be satisfied in a wide range of settings, and is certainly true for our running example given below. We also argue there that in

such scenarios $\bar{U}_1(d)$ will usually be *increasing* in the chosen intensity - so the higher the intensity of action the greater the expected quality of immediate success, and $\bar{U}_2(d)$ *decreasing* in d - so the more intense the action the greater the expected failure of the acceptability consequences of the action. The commander's adopted expected utility maximizing *Bayes decision* d^* will therefore trade off immediate success against acceptability score failure. He uses as much force to perform the mission task well whilst using as little force as possible to minimize jeopardizing the longer term issues, like retaining the support of the indigenous population.

Assume $\bar{U}_1(d)$ takes a value of 0 when $d < a_1$ and a value 1 when $d > b_1$ and $\bar{U}_2(d)$ takes a value of 0 when $d > a_2$ and a value 1, when $d < b_2$ where, by an abuse of notation, the lower bounds can be $-\infty$ and upper bounds ∞ . Then the commander believes that to use an intensity d below a_1 will surely have the worst immediate impact. However if d is chosen to be above b_1 he will surely be as successful in his immediate objectives as it is possibly to be. Similarly a_2 is the highest intensity he could use with minimal damage to the acceptability of his acts and b_2 is the lowest value at which the strategic objective would be most severely compromised. Henceforth call an action with intensity $d = b_1$ *pure combat* and an action $d = a_2$ *pure circumspection*. Pure combat is a decision a rational agent might choose were he to focus only on obtaining success in his immediate objectives, as measured through \bar{U}_1 . In contrast pure circumspection is a decision he might choose if he tried to obtain as high a score as possible on his acceptability objectives \bar{U}_2 .

To simplify the study of the geometry of $\bar{U}(d)$ in [4] we show that under the assumptions above $\bar{U}(d)$ can be conveniently written as a strictly increasing linear function of

$$V(d) = \exp(\rho)P_1(d) - P_2(d) \quad (2)$$

where

$$\begin{aligned} P_1(d) &\triangleq (u_1^+ - u_1^-)^{-1} (\bar{U}_1(d) - u_1^-) \\ P_2(d) &\triangleq 1 - (u_2^+ - u_2^-)^{-1} (\bar{U}_2(d) - u_2^-) \end{aligned} \quad (3)$$

and

$$\rho = \{\log k_1 - \log k_2\} + \{\log (u_1^+ - u_1^-) - \log (u_2^+ - u_2^-)\} \quad (4)$$

Because $V(d)$ is a location and scale change of $\bar{U}(d)$, note $\arg \max \bar{U}(d) = \arg \max V(d)$, so a commander's acts will be determined by the stationary points of $V(d)$.

To interpret the terms $\rho, P_1(d), P_2(d)$ in (2), first note from (3) that $P_1(d)$ and $P_2(d)$ are *distribution functions*, so each takes a minimum value 0 and maximum value 1. From the definitions above $P_i(d)$ has support $[a_i, b_i]$, $i = 1, 2$. The parameter ρ in (4) - here called the *daring* - takes values on the real line. The first term in (4) is increasing in k_1/k_2 . So the more weight the commander gives to the immediate success of his mission as compared to the acceptability of his action the larger the value of this parameter. The term $\log(u_1^+ - u_1^-)$ measures how much his choice of intensity can potentially impact on his immediate success whilst $\log(u_2^+ - u_2^-)$ measures the corresponding negative impact this intensity can have on the acceptability of his acts.

Because she is remote from the field of engagement, C2 may have difficulty accurately predicting how her commander assesses these latter impacts - and hence ρ - because, unlike him, she will not be fully aware of how the conflict is unfolding on the ground. However ρ impacts significantly on the decision a commander makes. As $\rho \rightarrow -\infty$ his expected utility will tend uniformly to $\bar{U}_2(d)$. Then pure circumspection a_2 tends to optimality. As ρ increases - for a given $(P_1(d), P_2(d))$ - so will the associated Bayes intensity $d^*(\rho)$ until as $\rho \rightarrow \infty$ pure combat tends to optimality [4].

A breakdown of rationality as well as contiguity challenges can arise whenever two close values of ρ can provoke the commander to choose two very different intensities of action. For then a commander's Bayes decision can jump up dramatically in response to an infinitesimal increase in ρ . A slight change in an agent's circumstances or confidence might then suddenly cause him to regret not committing to a very different decision. Furthermore in these circumstances two contiguous but autonomous commanders with similar expected marginal utilities, because they see similar information, may nevertheless act with very different intensities, causing one to retreat whilst the other attack. C2 can therefore explore how these scenarios she wishes to avoid might arise by examining the circumstances when $V(d)$ exhibits multiple maxima. It is surprising that by making some further additional fairly weak assumptions, it is possible to classify the geometry of $V(d)$ into 3 broad categories in terms of qualitative features that the remote C2 may have available to her [4]. The first category of scenarios are ones which can *inevitably* give rise to bipolar scenarios whilst the second describes when the commander will *never* be confronted with this problem. In this paper we will therefore focus on the last most interesting category where C2 can exert some control. This occurs when $a_1 \leq a_2 < b_1 \leq b_2$, when the commander can always find an SEU maximizing decision $d^* \in [a_2, b_1]$. Then, depending on how the C2 regulator communicates the mission, her commander

might choose an extreme decision close to pure circumspection or combat or alternatively some other intermediate intensity of engagement. We argue below that whenever possible a C2 regulator should try to construct scenarios where the commander will be predisposed to choose such an intermediate intensity.

3 An Example of a Military Conflict Decision

So consider the following example lying in this interesting category. Here $a_1 = a_2 = -1$ and $b_1 = b_2 = 1$ so without loss a commander can always be assumed to choose an intensity $-1 \leq d \leq 1$.

Example 1 *A battle group is set the immediate objective of securing two districts a, b of a city. Its commander believes that each district will take a minimum time to clear plus an exponentially distributed delay x_a and x_b with rate parameter β_1 . He believes he will have failed his immediate mission unless this task is fully completed. However before the securing a and b his strategic task is to first evacuate vulnerable civilians from two other districts c, d . There are potential delays x_c and x_d with rate parameter β_2 to add to the minimum time to complete this strategic objective and he believes he will have failed from a public acceptability viewpoint unless he is able to evacuate both areas successfully. The whole mission must not be delayed by more than 2 units of time. The commander must commit now to the time $d + 1$ he allows for delays in the immediate objective (so implicitly budgets for a $1 - d$ delay in completing the evacuation). He believes that $\beta_1 \perp \beta_2$ and that β_i has a gamma density*

$$\pi(\beta_i) = \frac{\lambda_i^{\alpha_i}}{\Gamma(\alpha_i)} \beta^{\alpha_i-1} \exp(-\lambda_i \beta), \quad > 0$$

where $\alpha_i, \lambda_i > 0$, so the delay he expects for each operation is $\mathbb{E}(\beta_i) = \alpha_i \lambda_i^{-1}$, $i = 1, 2$.

The sum of two independent exponential variables with the same rate β has a Gamma $G(2, \beta)$ distribution. Therefore his predictive density for the delay experienced in each of the tasks has a density $p_i(t)$, $t > 0$, with a unique mode $(\alpha_i + 1)^{-1} \lambda_i$, given by

$$p_i(t) = \int_{\beta>0} \beta^2 t \exp(-\beta t) \frac{\lambda_i^{\alpha_i}}{\Gamma(\alpha_i)} \beta^{\alpha_i-1} \exp(-\lambda_i \beta) d\beta = \frac{\Gamma(\alpha_i + 2) t \lambda_i^{\alpha_i}}{\Gamma(\alpha_i) [t + \lambda_i]^{\alpha_i+2}}$$

$i = 1, 2$. Since his utility function is zero-one on each attribute, his expected utility associated with using intensity d , $-1 \leq d \leq 1$, is proportional to

$$\exp(\rho) P_1(1 + d) + P_2(1 - d)$$

where $P_i(t)$ is the distribution function associated with $p(t)$, $i = 1, 2$. Any Bayes intensity must therefore satisfy

$$\exp(\rho) p_1(1+d) = p_2(1-d)$$

For illustration first suppose $\alpha_1 = \alpha_2 = 1$. and $\sigma \triangleq \exp(\rho) \lambda_1 \lambda_2^{-1} = 1$. Then this equation rearranges to the cubic

$$ad(1-d^2) + b(1-d^2) + cd + e = 0$$

where $a \triangleq 4 + 3(\lambda_1 + \lambda_2)$, $b \triangleq 3[\lambda_1(1 + \lambda_1) - \lambda_2(1 + \lambda_2)]$, $c \triangleq [\lambda_2^3 + \lambda_1^3]$ and $e = \lambda_1^3 - \lambda_2^3$. So letting $f \triangleq \frac{b}{3a}$, $g \triangleq (1 - a^{-1}c)$ and $z = d + f$ we obtain

$$z^3 + \{9f^2 + g\}z - f(g - 2f^2 - e) = 0$$

The local maxima of the commander's expected utility function can therefore be described by the well studied canonical cusp catastrophe [12], [9] where the splitting factor of this cusp catastrophe is $-(9f^2 + g)$. So if

$$9f^2 + g \geq 0 \iff b^2 \geq a(a - c)$$

the commander's expected utility function can have only one local maxima: a "intermediate" Bayes intensity. When parameters summarizing the commander's information and his values lie in this region his choice of optimal decision will be a smooth function of those parameters: similarly trained and tasked adjacent commanders will tend to act similarly. On the other hand if

$$9f^2 + g < 0 \iff \frac{b^2}{a^2} < 1 - \frac{c}{a} \iff b^2 < a(a - c)$$

then, for values of $f(g - 2f^2 - e)$ close to zero, his expected utility function will have two local maxima - one nearer minimal engagement ($d = -1$) the other maximal engagement ($d = 1$) with a minimum in between. Small changes in parameters may then cause different contiguous commanders to choose different decisions or a single commander to regret having committed to his chosen act.

So for example in the completely symmetric scenario when $\lambda_1 = \lambda_2 = \lambda$ there is always a stationary point at 0. When λ is very large (so that the expected delays are very small) 0 is the location of the unique stationary point. However whenever

$$2 + 3\lambda - \lambda^3 > 0$$

$d = 0$ corresponds to a *minimum* of the expected utility. It is easy to check that this inequality is satisfied only if $\lambda^{-1} > \frac{1}{2}$ i.e. when the sum

over the four expected delays is greater than the total time allowed for the mission. The larger this expected total delay is the more extreme the two equally preferable alternative decisions are. So the regulator should try to ensure that her commander is not faced with a potentially bipolar scenario where $9f^2 + g < 0$. In this symmetric scenario this inequality simply translates into C2 endeavoring to ensure the commander expects that he has been given enough time to complete both parts of his mission successfully.

4 The General Case

The example above, and the general principle it gives rise to, would be unremarkable if the geometry of the expected utility were heavily dependent on the distributional assumptions made by the commander. However it can be shown that *qualitatively* the geometry determining whether or not a given scenario has the potential to compromise rationality is surprisingly robust to changes in the algebraic forms of the commander's beliefs and values. So the sorts of broad criteria that a regulator should adopt, such as allowing the commander enough time to complete his mission successfully, can be determined in a much more general framework. In this general framework the conditions needed to avoid bipolar decision making are expressed in terms of points of inflexion of certain functions. Just as in the example above, the position of these points can in turn usually be expressed in a qualitatively meaningful way and be subject to the influence of the regulator.

Following [7], [10] we next investigate the geometrical conditions determining when dangers to loss of rationality might exist. Henceforth assume that the regulator believes that her commanders all have distributions P_i that are unimodal and twice differentiable in the open interval (a_i, b_i) , $i = 1, 2$ and constant nowhere in this interval. Any local maximum of $V(d)$ will then either lie on the boundary of the feasible space or satisfy

$$v(d) \triangleq f_2(d) - f_1(d) = \rho \quad (5)$$

where $f_i(d) = \log p_i(d)$, $i = 1, 2$. Provided the derivative $Dv(d) \geq 0$ such a stationary point must be a local maximum of V .

Let ξ_i denote the maximum (or mode) of $p_i(d)$ $i = 1, 2$. Because ξ_1 is a point of highest incremental gain in mission we call this point the *mission point* and the intensity ξ_2 where the threat to campaign objectives worsens fastest the *campaign point*. Note that it is not unreasonable to assume that two similarly trained and missioned commanders will entertain similar campaign and mission points when facing similar scenarios.

From this definition if $\xi_1 \leq \xi_2$, then we prove in [4] that for any $d \in [\xi_1, \xi_2]$, $v(d)$ is strictly decreasing. It follows that there is at most

one solution d^* to (5) for any value of ρ and $Dv(d) \geq 0$. Therefore this stationary value $d^* \in (a_2, b_1)$ is a local maximum of V . A regulator will find this scenario a manageable one: the mode of the intensity of immediate success is less than the modal intensity for acceptability success and there is a *unique* interior maximum in this interval: a rational agent will always choose an intermediate decision. So although their actions will depend on ρ , two commanders with similar but different daring ρ will act similarly. So if contiguous commanders are matched by their training and emotional history then it is plausible to conclude that they will make similar and hence broadly consistent choices.

On the other hand we prove in [4] that when $\xi_2 < \xi_1$ there is a value of ρ and an associated decision which is a stationary point of V for which the derivative of $v(d)$ is negative. This a stationary point will therefore be a local *minimum*. Then, just as in our example, the set of optimal decisions bifurcates into two disjoint sets: one lying in the interval of a "lower intensity" consistent with acceptability objectives, and the other in an interval of intensity favouring achieving the immediate mission objectives. At value of $\rho = \rho^\circ$ an option in each of these sets will be optimal giving rise to bipolarity in the decision space.

Example 2 *In our running example but when all prior parameters are arbitrary the condition on the two modes ξ_1, ξ_2 above implies bifurcation can occur for some value of ρ iff*

$$(\alpha_1 + 1)^{-1} \lambda_1 + (\alpha_2 + 1)^{-1} \lambda_2 < 2$$

The original interpretation that risks are avoided if the delays are perceived as not too large is therefore retained. When $\alpha_1 = \alpha_2 = \alpha$ and $\lambda_1 = \lambda_2 = \lambda$ this condition says that the prior expectation μ for dealing with any area satisfies $\mu \leq (1 + \alpha)^{-1}$ where $\alpha^{-1/2}$ is the coefficient of variation of the prior. So as $\alpha \rightarrow \infty$ there is less and less point aborting the evacuation because the commander believes he is unlikely to encounter any less delays in the pursuing the immediate elements of his mission.

From the geometrical arguments given above we can therefore argue that the phenomena we illustrated with our last example are generic and depend only on certain broad features such as that the predictive distributions associated with the two attributes being unimodal and smooth. The potential for the commander to act inappropriately simply depends on the relative positions of his campaign and mission points. Of course the C2 regulator may not be able to predict precisely when $\xi_1 > \xi_2$. However it will be possible to build a probabilistic model of this event using the information she has at hand or from general experiential knowledge. Such issues are discussed further in [4].

5 A Discussion of Further Elaborations

The running example above can also be generalized into a dynamic stochastic control problem where the commander has a choice of whether or not to abort the evacuation at *any* given time and simply focus on the immediate imperatives. Although the associated technicalities are beyond the scope of this paper, it appears that an optimal policy is straightforward to calculate. The qualitative behavior of the commander is analogous to the one above except that he tends to give up earlier on the evacuation if he is delayed longer than he predicted.

However a prior dependence between the parameters of the agent's two densities can affect C2's deductions significantly. For example if the agent assumes the delays are primarily due to a general lack of competence in his unit and not the situation on the ground then in the dynamic control setting above, instead of believing $\beta_1 \amalg \beta_2$ he might believe that $\beta_1 = \beta_2$. Were both his criterion weights $\frac{1}{2}$ then it is easily checked that pure circumspection is always optimal for him: i.e. he should continue evacuation until the two area are clear and only then address the immediate mission objective. Of course these alternative beliefs may well be predictable to a regulator who knows her commander.

If the commander does not have value independent attributes then this can also potentially giving rise to qualitatively very different behavior. For example he might believe his mission of securing areas c and d will have failed unless he has first fully completed the evacuation. In this case his utility function will be of the form $U = k_1 U_1 U_2 + k_2 U_2$ and it is easily checked that the decision to "continue until evacuation is finished and only then secure the other areas" is again always optimal. Note that such beliefs will again often be predictable to the regulator and indeed she can incline contiguity by communicating the mission in this way.

Finally note that the models we have developed above whose decisions lie on a continuum nevertheless induce a *finite discrete* reformulation capturing some of the qualitative features of the agent's appropriate behavior that might be predictable to his regulator. These finite states can in turn be used to define an implicit game between the two adversaries: each adversary attempting to induce irrational or discontinuous acts in the other: for example by deceiving the adversary into facing its commanders with bipolarity.

These various points illustrate that by making qualitatively different assumptions and by using this framework it is possible that she could come to quite different conclusions about how her agent will act. So even within the context as defined above, determining the appropriate ways to encourage obedient agent broadly to act effectively for C2 in any given scenario is a non trivial one requiring further detailed research.

On the other hand we have demonstrated that, provided the SEU assumption holds true for commanders not facing bipolar scenarios, C2 has a framework at least to investigate these issues. Furthermore the framework appears very robust in its conclusions to the misspecification of the factors needed to quantify fully her agent's assessments: factors which because of the devolution of decision making to her agent will be necessarily uncertain to the regulator.

References

- [1] Dodd, L. Moffat, J. Smith, J.Q. and Mathieson. G.(2003) "From simple prescriptive to complex descriptive models: an example from a recent command decision experiment" Proceedings of the 8th International Command and Control Research and Technology Symposium June Washington
- [2] Dodd, L.(1997), "Command decision studies for future conflict" DERA Unpublished Report.
- [3] Dodd, L. Moffat, J. and Smith, J.Q.(2006) "Discontinuities in decision-making when objective conflict: a military command decision case study" *J.Oper. Res. Soc.*,57, .643 - 654
- [4] Dodd, L. and Smith, J.Q. (2012) "Devolving Command Decisions in Complex Operations" *J.Oper. Res. Soc.*(to appear)
- [5] French, S. and Rios Insua, D.(2000) "Statistical Decision Theory" Arnold
- [6] Janis, J.L. and Mann, L.(1977) "Decision Making: A Psychological Analysis of Conflict, Choice and Commitment" Free Press. N.Y.
- [7] Moffat, J. (2002) "Command and Control in the Information Age" The Stationary Office, London
- [8] Moffat,J. and Witty, S. (2002) "Bayesian Decision Making and military control" *J. Oper. Res. Soc.* , 53, 709 - 718
- [9] Poston, T. and Stewart, I. (1978) "Catastrophe Theory and its applications" Pitman
- [10] Smith, J.Q., Harrison, P.J., and Zeeman, E.C.(1981) "The analysis of some discontinuous decision Processes" *E. J. Oper.Res.* Vol. 7, 30-43.
- [11] Smith, J.Q.(2010) "Bayesian Decision Analysis: Principles and Practice" Cambridge University Press
- [12] Zeeman E.C.(1977) "Catastrophe Theory: Selected Papers" Addison Wesley