

## A Extended context effect demonstrations

Here, we describe in detail how each of the context effects we mention in Section 4.1 can be elicited within the relative desirability framework. For ease of exposition, we reproduce the reference table described in the main paper as Table 2. In all the cases described below, we assume arbitrary context history sequences, but enforce equal observation counts for all contexts to match the standard epistemic assumptions of equal familiarity with all items. Further, relative desirability counts are notional, and it is possible in some cases to find partitions of relative desirability between options that render a particular demonstration invalid. We will point these out where relevant, and interpret them as theoretical substantiations of the observed fragility of these effects to changes in the relative desirabilities of the options in the initial choice set.

Effect name	Description	Assumptions
Frog legs	$c_1 \leftarrow \{X, Y\} \Rightarrow X \succ Y, c_2 \leftarrow \{X, Y, Z\} \Rightarrow Y \succ X$	-
Similarity	$c_1 \leftarrow \{X, Y\} \Rightarrow X \succ Y, c_2 \leftarrow \{X, Y, Z\} \Rightarrow Y \succ X$	$Z \approx X$
Attraction	$c_1 \leftarrow \{X, Y\} \Rightarrow X \sim Y, c_2 \leftarrow \{X, Y, Z\} \Rightarrow X \succ Y$	$X \succ Z$
Compromise	$c_1 \leftarrow \{X, Y\} \Rightarrow X \succ Y, c_2 \leftarrow \{X, Y, Z\} \Rightarrow Y \succ X$	$Y \succ^{(c)} X, Z$
Reference point	$c_1 \leftarrow \{X, Y\} \Rightarrow X \sim Y, c_2 \leftarrow \{X, Y, Z\} \Rightarrow X \succ^{(-)} Y$	$Z \succ X$

Table 2: A unified description of context effects.  $\succ$  indicates stochastic preference for one item over another.  $\succ^{(c)}$  indicates that the preference in question holds only in some observation contexts.  $\succ^{(-)}$  indicates that the preference in question is stochastically weaker than before.

### A.1 Similarity effect

In the similarity effect, given a choice set  $XY$ , a subject prefers option  $X$ . Now, a third option  $Z$  is introduced into the choice set, which is known to be similar to  $X$ , and not generally perceived to be clearly superior or inferior to  $X$ . In this expanded choice set  $XYZ$ , the subject is observed to reverse his preference and select  $Y$ . Let us say that the subject samples the choice set  $XY$  10 times in a row and prefers  $X$  6 times and  $Y$  4 times in those instances. At this point, from the agent’s standpoint, there is only one observable context  $XY$ , and hence, a simple derivation of the relative desirability follows  $R(X) = 0.6, R(Y) = 0.4$ , resulting in a rational preference for  $X$ . Now, the third option is introduced, with the agent possessing a history of previous comparisons between  $X$  and  $Z$ , but not  $Y$  and  $Z$ . We assume the subject’s indifference between  $X$  and  $Z$  to be reflected in desirability functions that consider either of these two options more desirable than the other half of the time in direct comparison. Thus, in 10 previous comparisons, we can estimate that each one will have been considered desirable 5 times. The set of observable contexts now expands to cover both  $XY$  and  $XZ$ . Let the posterior belief on the two contexts, dependent on the relative sequence of their past observations be  $p(c) = \{p, 1 - p\}$ . Finally, because of the similarity between  $X$  and  $Z$ , evidence for the desirability of  $X$  is computed over both observations of  $X$  and  $Z$  such that,

$$\begin{aligned}
 R(X) &= \frac{\sum_x^{\{X,Z\}} \sum_1^{10} p(r|x, XY)p(x|XY)p(XY) + \sum_x^{\{X,Z\}} \sum_1^{10} p(r|x, XZ)p(x|XZ)p(XZ)}{\sum_x^{\{X,Z\}} \sum_1^{10} p(x|XY)p(XY) + \sum_x^{\{X,Z\}} \sum_1^{10} p(x|XZ)p(XZ)}, \\
 &= \frac{6 \times 1 \times 1 \times p + 0 + 5 \times 1 \times 1 \times (1 - p) + 0}{10 \times 1 \times 0.5 + 10 \times 1 \times 0.5 + 10 \times 1 \times 0.5 + 10 \times 1 \times 0.5}, \\
 &= \mathbf{0.25 + 0.05p},
 \end{aligned}$$

while a similar computation for  $R(Y)$ ,

$$\begin{aligned}
 R(Y) &= \frac{\sum_1^{10} p(r|Y, XY)p(Y|XY)p(XY) + \sum_1^{10} p(r|Y, XZ)p(Y|XZ)p(XZ)}{\sum_1^{10} p(Y|XY)p(XY) + \sum_1^{10} p(Y|XZ)p(XZ)}, \\
 &= \frac{4 \times 1 \times 1 \times 0.5 + 0}{10 \times 1 \times 0.5 + 0}, \\
 &= \mathbf{0.4},
 \end{aligned}$$

Since  $p < 1, R(X) < 0.3 < R(Y)$ , resulting in a rational preference reversal  $Y \succ X$ . Recalculation using different values for the  $XY$  preference (e.g. 9/1 instead of 6/4) suggests that the

similarity effect will disappear in cases where  $X$  is clearly preferable to  $Y$ , an easily testable prediction from our theory. Further, the inference mechanism allows us to also predict that the similarity effect will return in such cases with the introduction of yet more items  $Z'$  similar to  $X$  into the choice set.

## A.2 Attraction effect

In the attraction effect, given a set of choices,  $\{X, Y\}$ , the subject is originally seen to be indifferent between the two options. However, when a third option  $Z$  that is similar to, but clearly inferior to option  $X$  is introduced into the choice set, the subject's preference switches to prefer option  $Y$ . As in the similarity effect, we interpret indifference between options to be reflected in desirability functions that consider either of these two options more desirable than the other half of the time. For 10 observations with the choice set  $\{X, Y\}$ , we will have 5 instances of  $X$  being more desirable than  $Y$  and 5 vice versa. As before, this leads to a simple relative desirability computation  $R(X) = 0.5, R(Y) = 0.5$ . Now, when we introduce  $Z$  into the choice set, and find that it is dominated by  $X$ , we obtain an additional supply of desirability observations from a second context  $XZ$  wherein  $X$  has almost always been found more preferable than  $Z$ . Let us say that, given 10 observations from this context, 8 favor  $X$  and 2 favor  $Z$ . Let the posterior belief on the two contexts, dependent on the relative sequence of their past observations be  $p(c) = \{p, 1 - p\}$ . Then, in the extended choice set regime, the desirability computation for  $X$  yields,

$$\begin{aligned} R(X) &= \frac{\sum_1^{10} p(r|X, XY)p(X|XY)p(XY) + \sum_1^{10} p(r|X, XZ)p(X|XZ)p(XZ)}{\sum_1^{10} p(X|XY)p(XY) + \sum_1^{10} p(X|XZ)p(XZ)}, \\ &= \frac{5 \times 1 \times 1 \times p + 8 \times 1 \times 1 \times (1 - p)}{10 \times 1 \times p + 10 \times 1 \times (1 - p)}, \\ &= \mathbf{0.8 - 0.3p}, \end{aligned}$$

while a similar computation for  $Y$  yields,

$$\begin{aligned} R(Y) &= \frac{\sum_1^{10} p(r|Y, XY)p(Y|XY)p(XY) + \sum_1^{10} p(r|Y, XZ)p(Y|XZ)p(XZ)}{\sum_1^{10} p(Y|XY)p(XY) + \sum_1^{10} p(Y|XZ)p(XZ)}, \\ &= \frac{5 \times 1 \times 1 \times 0.5 + 0}{10 \times 1 \times 0.5 + 0}, \\ &= \mathbf{0.5}. \end{aligned}$$

$p$  being a probability with non-zero support for both contexts  $p \in (0, 1) \Rightarrow R(X) > R(Y)$ , resulting in the establishment of a rational preference  $X \succ Y$  in place of the earlier indifference. This conclusion is expected to hold for any possible combinations of  $XZ$  preferences that clearly favor  $X$ .

## A.3 Reference point effect

The reference point effect has been used to explain many divergent sets of phenomena in the behavioral economics literature. For our demonstration, we restrict ourselves to explaining the results of a particular experiment on human subjects due to Vlaev et al [23], where subjects paid money to avoid forthcoming electric shocks of three different intensities, low, medium and high. The researchers found that subjects consistently paid more money to avoid pains that were greater than others in their recent history. In two sets of experiments, one where low shocks were mixed with medium shocks and one where medium shocks were mixed with higher ones, it was found that subjects paid much more money to buy out of medium shocks in the first condition than the second. Essentially, their pain evaluation was contingent on the set of pain options that they were being presented with. In the context-presentation framework, this can be posed as a problem where the subject is first offered the choice set  $LM$ , followed by further exposure to the choice set  $MH$ . Assuming that the subject experiences  $LM$  10 times, selecting  $M$  as the most painful option each time, their evaluation of relative (un)desirability  $R(M) = 1$ . Upon further presentation of 10 observations of  $MH$ , the new

desirability of  $M$  can be computed as,

$$\begin{aligned} R(M) &= \frac{\sum_1^{10} p(r|M, LM)p(M|LM)p(LM) + \sum_1^{10} p(r|M, MH)p(M|MH)p(MH)}{\sum_1^{10} p(M|LM)p(LM) + \sum_1^{10} p(M|MH)p(MH)}, \\ &= \frac{10p + 0}{10}, \\ &= p, \end{aligned}$$

which, being less than 1, implies that the (un)desirability of  $M$  reduces after exposure to a higher degree of pain. This observation, while almost trite on surface, has eluded the descriptive abilities of utility function approaches of measuring value, as described comprehensively in [22]. This effect is expected to remain stable for any choice of relative desirability frequency that respects the intuition that relatively lower levels of pain are more preferable.

#### A.4 Compromise effect

In the compromise effect, given a set of choices,  $\{X, Y\}$ , a subject prefers option  $X$ . Introduction of a third option  $Z$  leads to the development of two different ways of evaluating the desirability of any of the three items, resulting in situations where  $X$  may be strongly preferred to  $Z$  along one axis of measurement and strongly dominated by  $Z$  along the other. In standard descriptions of this effect, these different ways of evaluation are regarded as attributes, leading to a simple description of the problem in the framework of multi-attribute utility theory. For our current purpose, we achieve the same purpose notationally by considering  $XY$  and  $YX$  to be two different observation contexts representing possibilities that always co-occur, but are not always evaluated identically. Now, as in the earlier examples, at the time of the first observation, the only possible context is  $XY$ ; an observation history containing 7 preferences for  $X$  and 3 for  $Y$  results in a relative desirability calculation,  $R(X) = 0.7, R(Y) = 0.3$ .

Introduction of the third option, however, results in the (recalled) feasibility of six different contexts, which we index in  $\mathcal{C} = \{XY, YX, YZ, ZY, ZX, XZ\}$ . By the premise of the compromise effect setup, in the history of observing  $XZ$ ,  $X$  is preferred 8 times, while  $Z$  is preferred twice, while in observing  $ZX$ , these numbers are reversed. LSay observing 10 instances of  $YX$  yields 8 preferences for  $Y$  and 2 for  $X$ . 10 instances of  $YZ$  yield (6Y,4Z) while  $ZY$  yields (6Z,4Y). Since the contexts  $ij$  and  $ji$  are indistinguishable as observable choice sets, they occur with the same sample frequency. Thus, we can assume a posterior belief on six contexts,  $\{p_1, p_1, p_2, p_2, p_3, p_3\}$ . Then, upon observing  $XYZ$ , the desirability computation for  $X$  yields,

$$\begin{aligned} R(X) &= \frac{\sum_c^{\mathcal{C}} \sum_1^{10} p(r|X, c)p(X|c)p(c)}{\sum_c^{\mathcal{C}} \sum_1^{10} p(X|c)p(c)}, \\ &= \frac{7p_1 + 2p_1 + 0 + 0 + 2p_3 + 8p_3}{10p_1 + 10p_1 + 0 + 0 + 10p_3 + 10p_3}, \\ &= 0.05 \frac{9p_1 + 10p_3}{p_1 + p_3}. \end{aligned}$$

A similar computation for  $R(Y)$  yields,

$$\begin{aligned} R(Y) &= \frac{\sum_c^{\mathcal{C}} \sum_1^{10} p(r|Y, c)p(Y|c)p(c)}{\sum_c^{\mathcal{C}} \sum_1^{10} p(Y|c)p(c)}, \\ &= \frac{3p_1 + 8p_1 + 2p_2 + 8p_2 + 0 + 0}{10p_1 + 10p_1 + 10p_2 + 10p_2 + 0 + 0}, \\ &= 0.05 \frac{11p_1 + 10p_2}{p_1 + p_2}. \end{aligned}$$

Setting  $p_2 = p_3$  is equivalent to assuming  $Y$  and  $X$  both have equal histories of comparisons with the new option, which, while never a stated condition for observing the compromise effect, is not *prima facie* unreasonable. Doing so immediately forces  $R(X) < R(Y)$ , rendering the preference  $Y \succ X$  rational. The compromise effect has many more assumptions about relative preference frequencies than the attraction and similarity effect descriptions, rendering a comprehensive analysis intractable.

It is clear, however, that assuming a symmetric relationship between the XY and YX preferences, as we do in all the other cases, partially breaks the compromise effect, by rendering  $X \sim Y$ . Hence, we predict that a necessary requirement for the compromise effect to hold is for option Y to be more clearly preferable than option X along the new axis of evaluation introduced by inclusion of Z in the choice set.

## B Proof of representational equivalence

Here, we present a proof of the decision-equivalence between relative desirability and ordinal utility asserted in Section 4.2. Recall that proving this result is equivalent to proving that for any two possibilities  $x_i, x_j \in \mathcal{X}$ ,

$$x_i \succ x_j \Leftrightarrow R(x_i) > R(x_j). \quad (9)$$

Naturally, this will not be true in general for context-sensitive agents, since our framework specifically allows for preference reversals across multiple contexts, immediately rendering the LHS condition  $x \succ y$  insufficiently descriptive of preference relations. Therefore, we supplement it with a **context consistency** requirement,

$$\exists c \in \mathcal{C}, s.t. x_i \succ x_j \Rightarrow x_i \succ x_j \forall c \in \mathcal{C}_{ij}, \{x_i, x_j\} \in \mathcal{C}_{ij} \subseteq \mathcal{C}. \quad (10)$$

This additional requirement makes the expression of preferences in the context-aware setting epistemologically equivalent to the standard characterization of binary preference, since an observer insensitive to context will simply find that  $x_i \succ x_j$  whenever the two possibilities are observed together. To completely characterize a preference relation over  $\mathcal{X}$ , however, simply specifying consistent binary preferences is insufficient. Analogous to the regular concept of transitivity, we further assume the existence of **transitivity between contexts**, such that,

$$\text{if } x_i \succ x_j \text{ in } c_1 \text{ and } x_j \succ x_k \text{ in } c_2, \forall c \in \mathcal{C}, x_i \succ x_k, \quad (11)$$

thereby introducing a sense of preference order across observable contexts.

Now, consider that for any pair of possibilities  $\{x_i, x_j\} \subseteq \mathcal{X}$ , the set of observable contexts can be partitioned as,

$$\mathcal{C} = \mathcal{C}_{\setminus ij} \cup \mathcal{C}_{i \setminus j} \cup \mathcal{C}_{j \setminus i} \cup \mathcal{C}_{ij},$$

with the subscript indices indicating the possibilities from among  $\{x_i, x_j\}$  considered feasible, i.e.  $p(x|c) = 1$  within that context subset. Let  $\mathbb{C} = \{\mathcal{C}_{\setminus ij}, \mathcal{C}_{i \setminus j}, \mathcal{C}_{j \setminus i}, \mathcal{C}_{ij}\}$ . Then, we can expand the desirability definition in Equation (4) to,

$$R(x) = \frac{\sum_i |\mathbb{C}| \sum_c^{\mathbb{C}^{(i)}} p(r^{(t)}|x, c) p(x|c) p(c)}{\sum_i |\mathbb{C}| \sum_c^{\mathbb{C}^{(i)}} p(x|c) p(c)}, \quad (12)$$

Using our definitions of  $p(x|c)$  and  $p(r|x, c)$  (see (7) and immediately contiguous text), it is straightforward to show that,

$$R(x_i) = \frac{k_i \sum_c^{\mathcal{C}_{i \setminus j}} P(c) + k_{ij} \sum_c^{\mathcal{C}_{ij}} P(c)}{\sum_c^{\mathcal{C}_{i \setminus j}} P(c) + \sum_c^{\mathcal{C}_{ij}} P(c)}, \quad R(x_j) = \frac{k_j \sum_c^{\mathcal{C}_{j \setminus i}} P(c) + k_{ji} \sum_c^{\mathcal{C}_{ij}} P(c)}{\sum_c^{\mathcal{C}_{j \setminus i}} P(c) + \sum_c^{\mathcal{C}_{ij}} P(c)}, \quad (13)$$

since all other contributions disappear due to corresponding entries in  $p(x|c)$  being zero. Here, the single indexed  $k_i$  counts the number of times possibility  $x_i$  was considered the most desirable in contexts including  $x_i$  and excluding  $x_j$ ;  $k_j$  being defined symmetrically. The double-indexed  $k_{ij}$  counts the number of times  $x_i$  is considered the most desirable possibility in contexts where  $x_j$  is also believed to be present. Again,  $k_{ji}$  is defined symmetrically.

From (13) it should be clear that, in general, differences in the sampling of contexts in an agent's history of observations, measured, for instance, as variations in the size of the context subsets  $\mathbb{C}^{(i)}$  will render comparisons between desirability values undecidable<sup>4</sup>. Hence, to retain consistent preferences, we require an additional condition on the history of observation contexts that generate our

<sup>4</sup>To see why this must be the case, observe that for any two functions of homologous form to R such that  $\frac{\alpha k_i + k_{ij}}{\alpha + 1} = \frac{\beta k_j + k_{ji}}{\beta + 1} + \theta$ , with the  $k$  values fixed, it is always possible to find a new  $\beta' = \beta \left(1 + \frac{\theta}{k_j} + \frac{\theta}{k_j(\beta + 1)}\right) + 1$  that will reverse the inequality.

relative desirability measure. Specifically, we assume,

$$\forall x_i, x_j \in \mathcal{X}, \lim_{t \rightarrow \infty} |\mathcal{C}_{i \setminus j}| = |\mathcal{C}_{j \setminus i}|, \quad (14)$$

reflecting the intuition that there be no informative reason underlying the partial observability of world possibilities, i.e., partial observability occurs via random subset selection from  $\mathcal{X}$ . Note that this assumption, by symmetry, also implies

$$\lim_{t \rightarrow \infty} p(x|data^{(t)}) = U(x), \quad (15)$$

$U(\cdot)$  representing the uniform distribution.

Given this, in the infinite data limit, we obtain

$$\begin{aligned} p(x|data) &= \sum_c^{\mathcal{C}_{i \setminus j}} p(c) + \sum_c^{\mathcal{C}_{i \setminus j}} p(c) = \sum_c^{\mathcal{C}_{j \setminus i}} p(c) + \sum_c^{\mathcal{C}_{i \setminus j}} p(c) = U(x), \\ \Rightarrow \sum_c^{\mathcal{C}_{j \setminus i}} p(c) &= \sum_c^{\mathcal{C}_{i \setminus j}} p(c), \end{aligned}$$

obviating the necessity of further accounting for the denominators in (13).

It is now quite straightforward to demonstrate both directions of (8). First, assuming the left hand side of (8) immediately sets  $k_{ji} = 0$ . Further, using memorylessness,  $k_i$  can now be interpreted as determining the number of times  $x_i$  dominates all other possibilities in  $\mathcal{X} \setminus \{x_j\}$ ;  $k_j$  vice versa. By (11)  $x_i$  dominates all possibilities that  $x_j$  dominates, by (14) the number of observations over which either possibility can dominate is equal and by (15), in the limit of infinite decision samples, they will observe the same alternative possibilities, implying  $k_i \geq k_j$ . Since  $k_{ij} > 0^5$ , we directly have,

$$\begin{aligned} k_i \sum_c^{\mathcal{C}_{i \setminus j}} p(c) + k_{ij} \sum_c^{\mathcal{C}_{ij}} p(c) &\geq k_j \sum_c^{\mathcal{C}} p(c), \\ \Rightarrow R(x_i) &> R(x_j). \end{aligned}$$

Assuming the RHS of (9) to be true, adopting the selection rule  $\max_x R(x)$  proves the converse. Hence, contingent on the three assumptions we have specified above, the relative desirability based decision framework encodes relative preference relations equivalently well as ordinal utility functions.

---

<sup>5</sup> Assuming the LHS of (8) forces  $k_{ij}$  to be at least 1.