

Aversive learning and psychopathology

1 Associations between aversive learning processes  
2 and transdiagnostic psychiatric symptoms revealed  
3 by large-scale phenotyping

4 Toby Wise<sup>1,2,3</sup> & Raymond J Dolan<sup>1,2</sup>

5 <sup>1</sup>Wellcome Centre for Human Neuroimaging, University College London, London, UK

6 <sup>2</sup>Max Planck UCL Centre for Computational Psychiatry and Ageing Research, University College London,  
7 London, UK

8 <sup>3</sup>Division of the Humanities and Social Sciences, California Institute of Technology, Pasadena, CA, USA

9

10 Short title

11 Aversive learning and psychopathology

12 Contact

13 Toby Wise: [t.wise@ucl.ac.uk](mailto:t.wise@ucl.ac.uk); +44(0)203 1087538

14 Keywords

15 Anxiety, threat, computational modelling, learning, OCD, transdiagnostic

## Aversive learning and psychopathology

### 16 Abstract

### 17 Background

18 Symptom expression in a range of psychiatric conditions is linked to altered threat perception,  
19 manifesting particularly in uncertain environments. How precise computational mechanisms  
20 that support aversive learning, and uncertainty estimation, relate to the presence of specific  
21 psychiatric symptoms remains undetermined. 400 subjects completed an online game-based  
22 aversive learning task, requiring avoidance of negative outcomes, in conjunction with  
23 completing measures of common psychiatric symptoms. We used a probabilistic  
24 computational model to measure distinct processes involved in learning, in addition to inferred  
25 estimates of safety likelihood and uncertainty, and tested for associations between these  
26 variables and traditional psychiatric constructs alongside transdiagnostic dimensions. We used  
27 partial least squares regression to identify components of psychopathology grounded in both  
28 aversive learning behaviour and symptom self-report. We show that state anxiety and a  
29 transdiagnostic compulsivity-related factor are associated with enhanced learning from safety,  
30 and data-driven analysis indicated the presence of two separable components across our  
31 behavioural and questionnaire data: one linked enhanced safety learning and lower estimated  
32 uncertainty to physiological anxiety, compulsivity, and impulsivity; the other linked enhanced  
33 threat learning, and heightened uncertainty estimation, to symptoms of depression and social  
34 anxiety. Our findings implicate distinct aversive learning processes in the expression of  
35 psychiatric symptoms that transcend traditional diagnostic boundaries.

## Aversive learning and psychopathology

### 36 Introduction

37 Many core symptoms of mental illness are linked to learning about unpleasant events in our  
38 environment. In particular, symptoms of mood and anxiety disorders, such as apprehension,  
39 worry, and low mood can intuitively be related to altered perception of the likelihood of aversive  
40 outcomes. Indeed, the importance of altered threat perception is a feature of many diagnoses  
41 that extend beyond disorders of mood to encompass conditions such as psychosis<sup>1</sup> and eating  
42 disorders<sup>2</sup>. As a result, research into how individuals learn about aversive events holds great  
43 promise for enhancing our understanding across a diverse range of mental health problems.

44 Computational approaches are a powerful means to characterise the precise mechanisms  
45 underpinning learning, as well as uncovering how these relate to psychiatric symptom  
46 expression<sup>3,4</sup>. Recent studies have leveraged computational modelling to capture associations  
47 between learning processes and psychiatrically-relevant dimensions in non-clinical samples<sup>5-</sup>  
48 <sup>8</sup>, as well as in clinical conditions ranging from anxiety and depression to psychosis<sup>9-12</sup>. A  
49 common finding across studies is that of altered learning rates, where psychopathology is  
50 linked to inappropriate weighting of evidence when updating value estimates<sup>7,13,14</sup>. Notably,  
51 there is evidence suggesting that people with clinically significant symptoms of anxiety and  
52 depression show biased learning as a function of the valence of information, updating faster in  
53 response to negative than positive outcomes presented as monetary losses and gains<sup>12</sup>, a bias  
54 that might engender a negative view of the environment. However, we previously found an  
55 opposite pattern in a non-clinical study using mild electric shocks as aversive stimuli, whereby  
56 more anxious individuals learned faster from safety than from punishment, and underestimated  
57 the likelihood of aversive outcomes<sup>15</sup>. This latter finding highlights a need for a more extensive  
58 investigation using larger samples.

59 In addition to aberrant learning another process implicated in the genesis of psychiatric disorder  
60 relates to the estimation of uncertainty<sup>16</sup>. While there are multiple types of uncertainty, here we  
61 use the term to refer to estimation uncertainty, describing the precision of a learned association.  
62 Estimation uncertainty is highest when there is a lack of experience, or the association to be  
63 learned is unstable. For example, having seen two coin flips and observing one head and one  
64 tail, one might believe the likelihood of observing a head is 50%, though you are highly  
65 uncertain about this estimate due to a lack of evidence. This kind of uncertainty plays a  
66 fundamental role in learning, and computational formulations optimise learning in the face of  
67 non-stationary probabilistic outcomes based on uncertainty<sup>11,17-20</sup>. While psychiatric  
68 symptoms, including anxiety, have been linked to an inability to adapt learning in response to  
69 environmental statistics such as volatility<sup>5,9</sup>, little research has investigated how individuals  
70 estimate, or respond to, uncertainty in aversive environments and its potential association with  
71 psychiatric symptoms. This is a crucial question given that core features of anxiety revolve  
72 around a concept of uncertainty. For example, individuals with anxiety disorders report feeling  
73 more uncertain about threat and being less comfortable in situations involving uncertainty<sup>21-24</sup>.  
74 Surprisingly, in an earlier lab-based study we observed a surprising relationship, finding that  
75 more anxious individuals were more certain about stimulus-outcome relationships<sup>15</sup>. However,  
76 this was in a relatively small sample and therefore warrants further investigation.

## Aversive learning and psychopathology

77 Existing work on aversive learning has had a particular focus on symptoms of anxiety and  
78 depression<sup>7,12</sup>. However, these approaches have not been designed optimally for identifying  
79 mechanisms that span traditional diagnostic boundaries. This assumes importance in light of  
80 recent studies, using large samples, showing several aspects of learning and decision making  
81 relate more strongly to transdiagnostic factors (symptom dimensions that are not unique to any  
82 one disorder) than to any specific categorical conception of psychiatric disorder<sup>6,8,25–27</sup>. Applying  
83 such an approach to aversive learning may yield better insights into the role of learning in  
84 psychiatric disorders. Additionally, computationally-defined measures of learning and decision  
85 making can facilitate identification of novel transdiagnostic factors, going beyond those  
86 identified based solely on correlated symptom clusters in self-report and clinical interview  
87 measures<sup>6,28–30</sup>.

88 Here, we aimed to clarify the nature of the relationship between aversive learning processes  
89 and traditional measures of anxiety, as well as transdiagnostic psychiatric factors identified in  
90 prior work<sup>6</sup> in a large, preregistered study conducted online. This allowed us to measure effects  
91 with high precision, potentially helping to resolve mixed findings from previous studies<sup>12,15</sup>, in  
92 addition to identifying small but meaningful effects that cross traditional diagnostic  
93 boundaries<sup>6</sup>. Thus, we used a computational approach to test whether anxiety and  
94 transdiagnostic symptoms are associated with biased learning from safety and threat, whether  
95 these factors relate to altered estimates of threat likelihood, and whether they are associated  
96 with different levels of uncertainty during threat learning. We then used partial least squares  
97 (PLS) regression, a data-driven multivariate method, to derive transdiagnostic latent  
98 components of psychopathology grounded in both self-report and computational measures.  
99 Given difficulties in using traditional aversive stimuli in an online setting, we developed a game-  
100 based avoidance task designed to engage threat and avoidance processes without the need  
101 for administration of painful or noxious stimuli. Both the task and modelling are, in principle,  
102 similar to our previous lab-based task<sup>15</sup>, but their implementation here allows straightforward  
103 administration in large samples recruited online.

104

## 105 Results

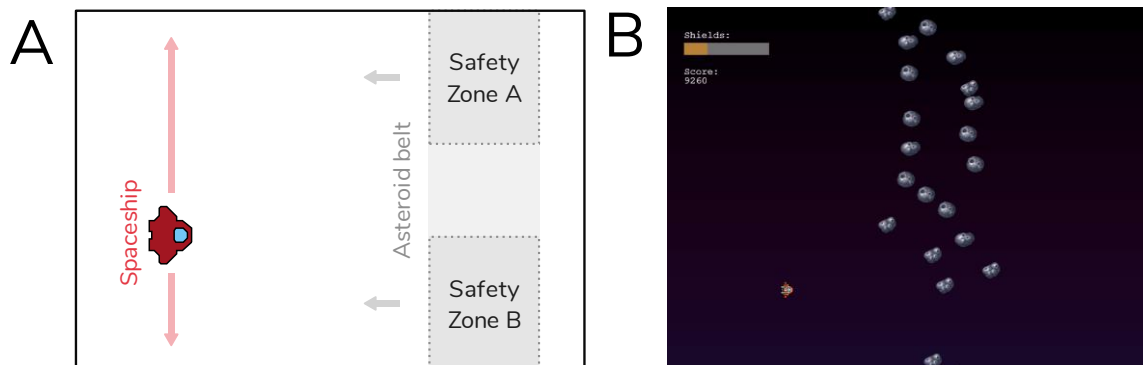
### 106 Task performance

107 Four hundred subjects recruited online through Prolific<sup>31</sup> performed a game-based aversive  
108 learning task, where the aim was to fly a spaceship through asteroid belts without being hit  
109 (Figure 1). Getting hit by the asteroids reduced the integrity of the spaceship, and after  
110 sufficient hits the game terminated. Crucially, there were two zones at the top and bottom of  
111 the screen where subjects could encounter a hole in the asteroid belt, each associated with a  
112 changing probability of being safe. In order to perform well at the task subjects needed to learn  
113 which zone was safest and behave accordingly.

114 Subjects were engaged and performed well at the task, with a median number of spaceship  
115 destructions of 1 (Interquartile range = 2) over the course of the task. They also reported high

## Aversive learning and psychopathology

116 motivation to perform the task, providing a mean rating of 85.70 (SD = 18.44) when asked to  
117 rate how motivated they were to avoid asteroids on a scale from 0-100.



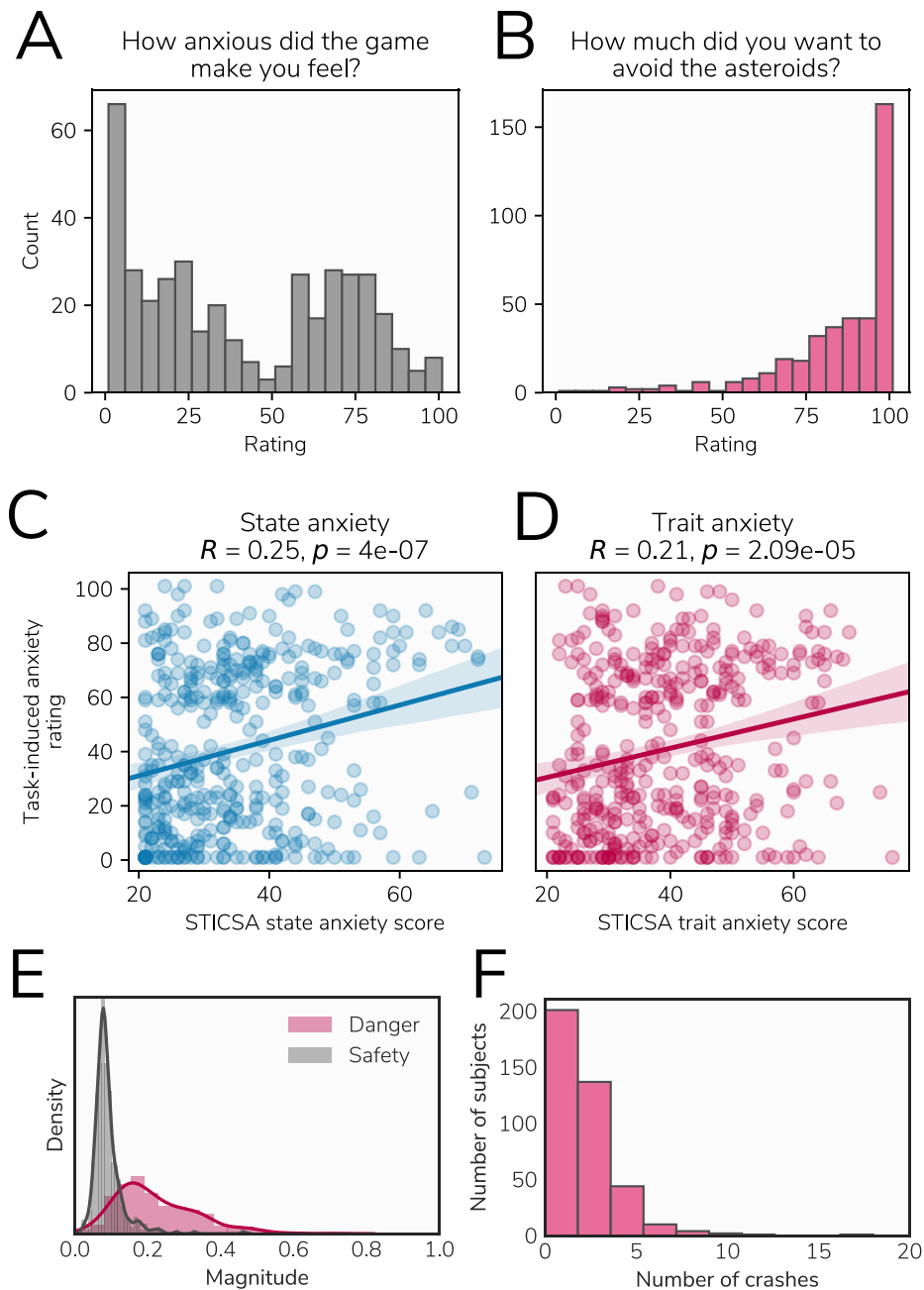
118

119 Figure 1. A) Task design. Subjects were tasked with playing a game that had a cover story involving flying a  
120 spaceship through asteroid belts. Each asteroid belt featured two locations that could potentially contain escape  
121 holes (safety zones), and subjects were instructed to aim to fly their spaceship through these to gain the highest  
122 number of points. Subjects were only able to move the spaceship in the Y-dimension, while asteroid belts moved  
123 towards the spaceship. The probability of each zone being safe varied over the course of the task but this could be  
124 learned, and learning this probability facilitated performance. B) Screenshot of the task, showing the spaceship, an  
125 asteroid belt with a hole in the lower safety zone (safety zone B), a representation of the spaceship's integrity (shown  
126 by the coloured bar in the top left corner) and the current score.

### 127 Computational modelling of behaviour

128 To quantitatively describe behaviour, we fit a series of computational models to subjects'  
129 position data during the task (see Methods and Supplementary Material for a full description of  
130 tested models). The winning model was a probabilistic model incorporating different updates  
131 parameters for safety and danger, as well as a "stickiness" parameter representing a tendency  
132 for subjects to stick with their previous position. This model represents an extension of one we  
133 have previously used successfully as a lab-based aversive learning task<sup>15</sup>, and is described fully  
134 in the Methods section. Briefly, this winning model assumes that subjects in the task represent  
135 the safety probability of each zone using a beta distribution, which is updated on each trial  
136 based on encounters with danger or safety. Simulating responses using the model, using each  
137 subject's estimated parameter values, produced behavioural profiles that demonstrated a high  
138 concordance with the true data, reproducing broad behavioural patterns seen in the true data  
139 (Figure S1).

## Aversive learning and psychopathology



140  
141 Figure 2. A) Distribution of task-induced anxiety ratings recorded after the task. B) Distribution of task motivation  
142 ratings. C and D) Relationships between task-induced anxiety ratings and state and trait anxiety scores. E) Degree  
143 of location switching after encountering danger and safety across subjects. The switch magnitude is the average  
144 absolute change in position between trial  $n$  and trial  $n+1$ . As expected, subjects showed more switching behaviour  
145 after encountering danger and were more likely to stay in the same position following a safe outcome. F) Distribution  
146 of crash number (representing the number of subjects hit enough asteroids to end the game) across subjects.

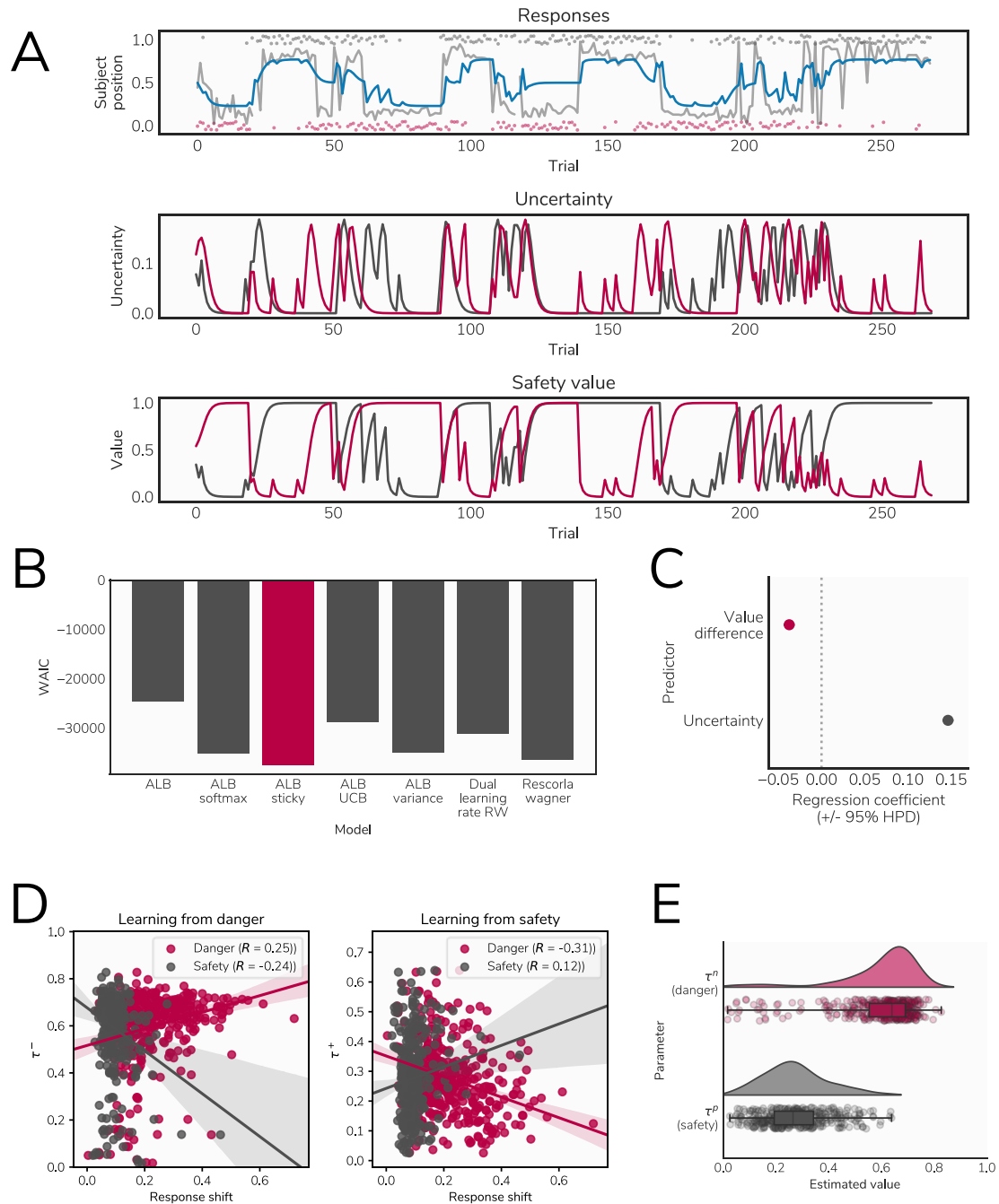
### 147 Task and model validation

148 In the light of the task's novelty, it was important to ensure the task has content validity, and  
149 that it produces behaviour reminiscent of more traditional tasks. Likewise, the computational  
150 models used should provide measures and parameter estimates that reflect the behaviour they  
151 aim to describe. We therefore conducted extensive validation exercises. These are reported  
152 fully in supplementary material, but we summarise these here.

## Aversive learning and psychopathology

153 First, we ensured the task did induce states of subjective anxiety in the majority of subjects  
154 (Figure 2A), and this level of anxiety was correlated with self-report state and trait anxiety  
155 (Figure 2C and 2D). Importantly, the task produced behaviour reminiscent of more traditional  
156 lab-based tasks, with subjects adjusting their position to a greater extent following danger than  
157 following safety (Figure 2E), as demonstrated in prior studies<sup>12,15,32</sup>. With respect to our  
158 computational model, we verified the model's update parameters were robustly correlated with  
159 subjects' tendency to move, or stay, following danger and safety respectively (Figure 3D). We  
160 also ensured that the safety value and uncertainty values produced by simulating data from  
161 our model with best fitting parameters correlated with subjects' model-free behaviour, finding  
162 that subjects changed their position more when model-derived uncertainty was high, and  
163 when the difference between the safety value of the two zones was small, as expected (Figure  
164 3C). Finally, we verified that our model's update parameters showed greater updating from  
165 danger relative to safety, as we found in a previous lab-based study<sup>15</sup>, finding this was indeed  
166 the case (Figure 3E).

## Aversive learning and psychopathology



167

168 Figure 3. A) Data generated from the model. The top panel shows responses and model fit for an example subject.  
 169 In the top panel, the grey line represents the subject's position throughout the task, with the grey and red dots  
 170 representing safe locations on each trial. The blue line represents simulated data from the model for this subject.  
 171 The lower two panels show estimated uncertainty and safety probability for each stimulus (represented by the grey  
 172 and red lines) across the duration of the task, generated by simulating data from the model. B) Model comparison  
 173 results, showing the WAIC score for each model with the winning model highlighted. ALB = asymmetric leaky beta,  
 174 RW = Rescorla-Wagner. C) Results of our analysis validating the safety value and uncertainty measures, showing  
 175 the extent to which each measure predicted subjects' tendency to switch position (described in supplementary  
 176 materials). D) Correlations between estimated update parameters for danger (left) and safety (right) and our  
 177 behavioural measure of position switching after these outcomes across subjects, demonstrating that parameters  
 178 from our model reflect purely behavioural characteristics. E) Distributions of estimated parameter values for  $\tau^+$  and



## Aversive learning and psychopathology

179  $\tau^-$ , representing update rates following danger and safety outcomes respectively, showing a bias in updating  
180 whereby subjects update to a greater extent in response to danger than safety.

181

### 182 Relationships with psychiatric measures

183 First, we asked whether our four behavioural variables of interest (threat update parameter,  
184 danger update parameter, mean estimated safety probability, and mean estimated uncertainty)  
185 were associated with anxiety (both state and trait) and intolerance of uncertainty. The strongest  
186 relationships, with HPD intervals that did not include zero, were positive effects of state anxiety  
187 on safety update rates and mean estimated safety probability (Figure 4, Table 1), although  
188 effects for trait anxiety were in the same direction and of a similar magnitude for some  
189 measures, indicating more anxious individuals learned faster about safety and perceived safety  
190 as more likely overall.

Target variable	Predictor	Estimate (+/- 95% HPDI)
Threat update ( $\tau^-$ )	IUS	0.07 (-0.03, 0.16)
	STICSA S	-0.09 (-0.19, 0.01)
	STICSA T	-0.03 (-0.12, 0.07)
Safety update ( $\tau^+$ )	IUS	0.0 (-0.1, 0.09)
	STICSA S	<b>0.12 (0.02, 0.21)</b>
	STICSA T	0.09 (-0.01, 0.18)
Mean safety uncertainty	IUS	-0.05 (-0.14, 0.05)
	STICSA S	-0.09 (-0.19, 0.01)
	STICSA T	-0.08 (-0.18, 0.01)
Mean safety probability	IUS	-0.02 (-0.12, 0.07)
	STICSA S	<b>0.12 (0.02, 0.21)</b>
	STICSA T	0.06 (-0.03, 0.15)

191 Table 1. Estimates from regression model predicting learning-related variables derived from our computational  
192 model from measures of intolerance of uncertainty, state anxiety, and trait anxiety. Effects with HPDIs excluding  
193 zero are shown in bold. IUS: Intolerance of uncertainty scale; STICSA S: State-trait inventory of cognitive and  
194 somatic anxiety, state measure; STICSA T: State-trait inventory of cognitive and somatic anxiety, trait measure;  
195 HPDI: Highest posterior density estimate

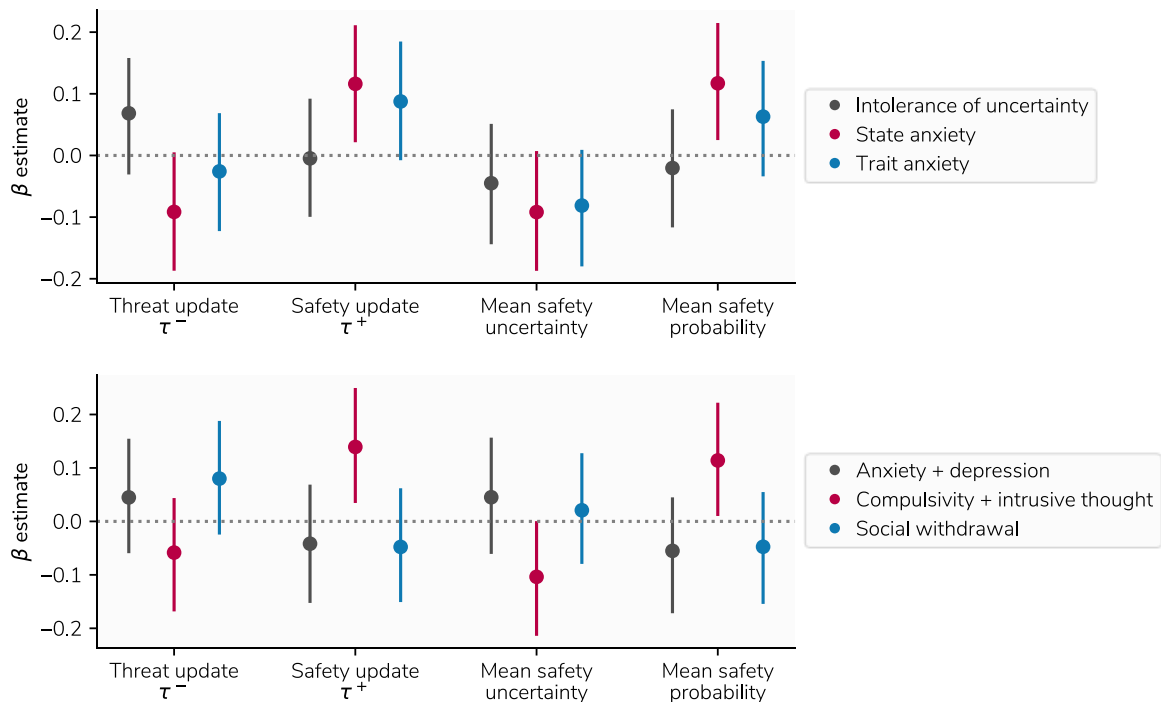
196 We then examined the extent to which task behaviour was associated with three  
197 transdiagnostic factors of psychopathology identified through self-report assessments in  
198 previous research<sup>6</sup>. Here, we observed effects of a factor labelled compulsivity and intrusive  
199 thought (Figure 4, Table 2), reflecting the fact that subjects scoring higher on this factor learned  
200 faster about safety and had higher safety probability estimates. There was also a weak effect  
201 of this factor on uncertainty, although the HPDI for this included zero. Other effects were weak,  
202 and including reported task motivation as a covariate had a negligible effect on the results (see  
203 supplementary results). Importantly, all of these analyses were determined *a priori* and are  
204 included in our preregistration.

205

## Aversive learning and psychopathology

Target variable	Predictor	Estimate (+/- 95% HPDI)
Threat update ( $\tau^-$ )	AD	0.04 (-0.06, 0.15)
	CBIT	-0.06 (-0.17, 0.04)
	SW	0.08 (-0.02, 0.19)
Safety update ( $\tau^+$ )	AD	-0.04 (-0.15, 0.07)
	CBIT	<b>0.14 (0.03, 0.25)</b>
	SW	-0.05 (-0.15, 0.06)
Mean safety uncertainty	AD	0.05 (-0.06, 0.16)
	CBIT	-0.1 (-0.21, 0.0)
	SW	0.02 (-0.08, 0.13)
Mean safety probability	AD	-0.06 (-0.17, 0.04)
	CBIT	<b>0.11 (0.01, 0.22)</b>
	SW	-0.05 (-0.15, 0.05)

206 Table 2. Estimates from regression model predicting learning-related variables derived from our computational  
 207 model from the three transdiagnostic factors identified by Gillan et al. (2016). Effects with HPDIs excluding zero  
 208 are highlighted. AD: Anxious-depression; CBIT: Compulsive behaviour and intrusive thought; SW: Social  
 209 withdrawal; HPDI: Highest posterior density estimate



210 Figure 4. Top panel: Results of state/trait anxiety and intolerance of uncertainty models, showing relationships  
 211 between these psychiatric variables and behavioural variables. Points indicate the mean of the posterior distribution  
 212 for the regression coefficient parameter, while error bars represent the 95% highest posterior density interval. The  
 213  $\beta$  estimate here refers to the regression coefficient for each predictor. Bottom panel: Results of three factor model,  
 214 showing relationships between behaviour and factors labelled anxiety and depression, compulsivity and intrusive  
 215 thought, and social withdrawal.  
 216

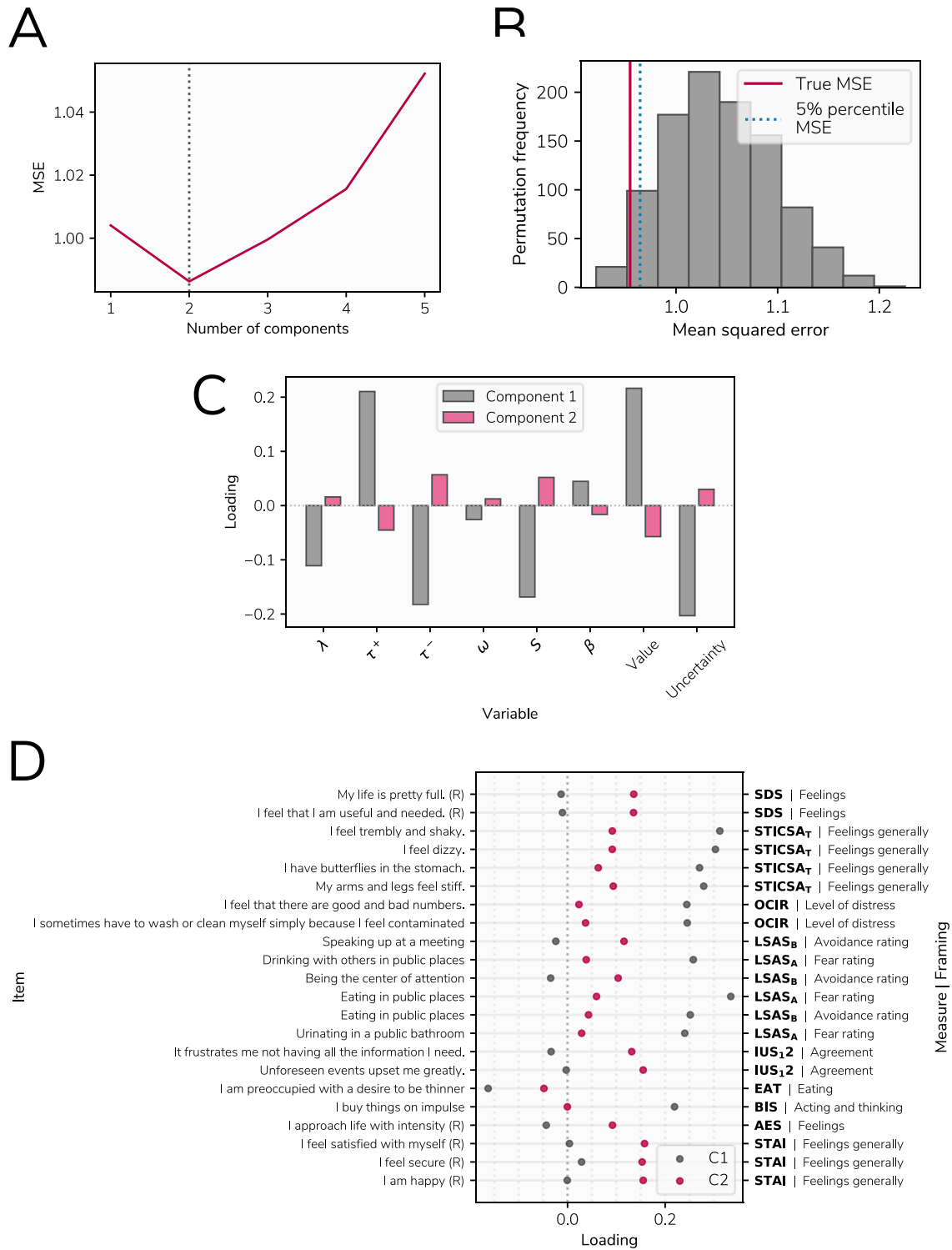
## Aversive learning and psychopathology

### 217 Psychiatric constructs derived from behaviour and self-report

218 Numerous studies have used dimensionality reduction procedures such as factor analysis on  
219 questionnaire-based data to identify factors of psychopathology that cut across diagnostic  
220 boundaries<sup>6,28–30</sup>. This, in turn, has revealed that many behaviourally-defined phenotypes are  
221 more strongly associated with transdiagnostic factors than any single disorder<sup>6,8,26</sup>. We built  
222 upon this work by incorporating computationally-derived indexes of behaviour into this  
223 dimensionality reduction procedure, where the aim was to identify latent constructs grounded  
224 in both self-report and behaviour. We used partial least squares (PLS) regression, a method  
225 that identifies latent components linking multivariate data from multiple domains based on their  
226 shared covariance. This method has been employed successfully to provide insight into how  
227 panels of cognitive and behavioural measures relate to multivariate neuroimaging-derived  
228 phenotypes<sup>33–35</sup>. We first identified the number of components that best describe our data by  
229 evaluating the performance of a predictive PLS model using cross-validation. We found two  
230 latent components gave the best predictive performance (Figure 5A). We then evaluated the  
231 performance of this model on held out data using permutation testing, showing our model  
232 achieved a statistically significant level of predictive accuracy (permutation  $p = 0.025$ , Figure  
233 5B). This indicates that our combined self-report and behavioural data is best explained by a  
234 two-component structure linking these two domains. Importantly, the fact that this level of  
235 accuracy was found on unseen data ensures that our results do not result from overfitting the  
236 training data<sup>36</sup>.

237 To aid interpretation of these two components we examined how behavioural variables loaded  
238 on each. The first component had positive weights on update rates in response to safety and  
239 estimated safety likelihood, and negative weights on update rates in response to threat, decay,  
240 stickiness, and mean uncertainty estimates (Figure 5C), while the reverse was true of the  
241 second component. Loadings on questionnaire items were varied, and labelling such  
242 components is invariably subjective. Nevertheless, the first component tended to load most  
243 strongly on items describing physical symptoms of anxiety, compulsive behaviour, and  
244 impulsivity. In contrast, the second latent component loaded primary on items describing social  
245 anxiety and depressed mood. For illustrative purposes, items with the top 10% percent of  
246 differences in loadings between components are shown in Figure 5D, with full details available  
247 in supplementary material.

## Aversive learning and psychopathology



248

249 Figure 5. Results of PLS regression analysis. A) The optimal number of components was determined based on cross-  
 250 validated predictive accuracy within 75% of the data used for training the model. This figure represents the negative  
 251 mean squared error of these predictions across models with between one and five factors, showing best  
 252 performance with 2 components. B) Null distribution of predictive accuracy scores generated by retraining our PLS  
 253 regression model on 1000 permuted datasets and testing on the held out 25% of the data set, with the mean  
 254 squared error (MSE) achieved by the model trained on the true data shown by the red line. C) Loadings for the two  
 255 components on behavioural variables, including all parameters in the model and mean safety probably and  
 256 uncertainty estimates, across all subjects. D) Loadings on questionnaire items showing the largest dissociations in

## Aversive learning and psychopathology

257 loadings between the two components, identified by taking the lowest and highest 10% of differences between  
258 loadings. Items marked (R) are reverse coded. The labels on the right indicate the measure the item is taken from  
259 and an indicator of how the question is framed. SDS: Zhung Self-Rating Depression Scale, STICSA: State Trait  
260 Inventory of Cognitive and Somatic Anxiety, OCIR: Obsessive Compulsive Inventory, LSAS: Liebowitz Social Anxiety  
261 Scale (A and B represent subscales), IUS12: Intolerance of Uncertainty Scale, EAT: Eating Attitudes Test, AES:  
262 Apathy Evaluation Scale, STAI: State Trait Anxiety Inventory.

## 263 Discussion

264 Perceptions of danger and safety have been linked to key symptoms of psychiatric disorders.  
265 Here, in a large-scale study examining aversive learning we show that when subjects learn to  
266 avoiding threat, transdiagnostic components of psychopathology relate to how they learn  
267 about both safety likelihood, and uncertainty.

268 We found a counter-intuitive relationship between biases in learning and the presence of  
269 features of anxiety. Subjects scoring higher on state anxiety tended to update their predictions  
270 to a greater extent in response to safety, as well as perceiving safety to be more likely overall,  
271 than those scoring low on this measure. These results diverge from previous findings that  
272 report individuals diagnosed with clinical anxiety and depression learn faster from  
273 punishment<sup>12</sup>, but are in concordance with our previous work in a non-clinical sample using a  
274 more traditional lab-based aversive learning task<sup>15</sup>. The large sample size employed here  
275 allowed us to estimate these effects precisely, making it unlikely that they are simply a product  
276 of statistical noise. One explanation for the discrepancy between our results and those found  
277 by Aylward et al.<sup>12</sup> is that this previous study included subjects with a mix of anxiety and  
278 depressive disorders, and a negative bias in learning may be more characteristic of depressive  
279 symptoms. Our PLS analysis provides some support for this speculation as we found that  
280 symptoms of depression were associated with elevated learning from threat, suggesting that  
281 such a bias in learning is associated more with depressive symptoms. It is also possible that  
282 the nature of our game-based online task engaged processes distinct from that of standard  
283 lab-based tasks. However, we believe this is unlikely since we replicate behavioural patterns  
284 shown in more traditional tasks, and also observed similar associations with anxiety in a  
285 previous lab-based study<sup>15</sup>. As such, we are confident that this is not simply due to the task  
286 used.

287 We found a similar pattern of enhanced learning from safety when examining a transdiagnostic  
288 factor representing compulsivity and intrusive thought. Although this factor has been shown  
289 to be associated with less model-based behaviour<sup>6,25</sup>, altered confidence judgements<sup>26</sup>, and  
290 action-confidence coupling<sup>8</sup> in large-scale samples, to date it has not been investigated with  
291 regard to threat learning. Notably, we also found a weak relationship between this factor and  
292 uncertainty, whereby more compulsive individuals had higher certainty in their safety  
293 estimates, echoing previous work in perceptual decision making that showed this factor is  
294 associated with higher confidence estimates<sup>8,26</sup>. We only found weak relationships (where the  
295 posterior density estimate crossed zero) with the other two factors, representing anxious-  
296 depression and social withdrawal, in a direction indicative of lower safety probability estimates  
297 and higher uncertainty.

## Aversive learning and psychopathology

298 Perhaps our most striking results come from a data-driven approach where we derived  
299 components of psychopathology grounded in computational analyses of aversive learning  
300 behaviour. This method provides conceptually similar results to the factor analytic methods  
301 used in previous large-scale online studies<sup>6</sup>, but builds upon this work by incorporating  
302 behaviour into the process. Using PLS regression, we identified two latent components, one  
303 broadly associating greater learning from safety with physiological symptoms of anxiety and  
304 compulsivity, while the other associated greater learning from threat with depressive  
305 symptoms and social anxiety. Notably, this data-driven analysis also revealed relationships  
306 between aversive learning and impulsive behaviour, encompassing a symptom dimension that  
307 is typically studied in the context of reward processing<sup>37</sup>. Individuals scoring higher on these  
308 symptoms exhibited higher safety learning, which may explain previously observed  
309 relationships between impulsivity and risk tolerance<sup>38</sup>. Importantly, while this analysis was  
310 exploratory, we demonstrate its robustness through testing on held-out data, ensuring our  
311 results are not affected by overfitting<sup>36</sup>.

312 Overall, the present results add to the growing literature showing associations between  
313 psychopathology and learning under uncertainty. Previous studies using computational  
314 approaches have largely focused on learning about rewards and losses<sup>10-12,27,39</sup>, or perceptual  
315 learning<sup>9</sup>, and those that have used more aversive paradigms (using outcomes intended to  
316 evoke subjective anxiety), such as learning to predict electric shocks, have been limited by small  
317 samples<sup>5,15,18,40</sup>. As a result, the precise role played by aversive learning processes in psychiatric  
318 symptoms has been unclear. Our work adds to this literature by providing an account of how  
319 these processes relate to symptoms across a range of traditional diagnostic categories.

320 The results we report raise questions of importance for future work. In particular, the finding  
321 that more anxious individuals tend to overestimate safety likelihood runs counter to intuition,  
322 and further work is required to understand how this may relate to symptom expression. One  
323 speculative possibility is that a persistent underestimation of threat likelihood would lead to an  
324 abundance of aversive prediction errors, causing a state of subjective anxiety. An alternative  
325 explanation is the result reflects a tendency for highly anxious individuals to seek safety, and  
326 be resistant to leaving places associated with safety<sup>41,42</sup>. However, these hypotheses await  
327 direct testing, and it will be especially important to examine them in large-scale clinical samples,  
328 taking into account a broader range of psychiatric phenotypes.

329 A further important feature of this study is our development of a new online task for measuring  
330 aversive learning. A number of studies examining other aspects of learning and decision  
331 making in the context of psychiatric disorders have also availed of large samples recruited  
332 through online services<sup>6,8,25,26</sup>. However, it has been difficult to examine aversive learning in  
333 online environments, as aversive lab stimuli such as shock cannot be easily administered online.  
334 Only one study thus far has investigated threat-related decision making (although not learning)  
335 online, using monetary loss as an aversive stimulus<sup>27</sup>. A game-based design allowed us to  
336 design a task that required avoidance behaviour as well as evoke feelings of anxiety, taking  
337 advantage of the well-known ability of games to produce strong emotional reactions<sup>43-47</sup>,  
338 resulting in a paradigm which we believe provides a more valid assessment of aversive learning

## Aversive learning and psychopathology

339 than more commonly used monetary-loss based tasks. Although qualitatively different from  
340 standard lab-based tasks, we observed similar patterns of biased learning to that seen in lab-  
341 based work<sup>15</sup>. An added benefit of our task is that it is highly engaging, and subjects reported  
342 feeling motivated to perform well. These features are not only important for the kind of large-  
343 scale online testing performed here. This task renders it feasible to measure aversive learning  
344 at regular intervals without subjects needing to physically visit the lab, a feature that could be  
345 of considerable utility in clinical trials.

346 One potential limitation of this study is a focus on a general population sample which, being  
347 recruited online, was not subject to the kind of detailed assessment possible offline. While this  
348 might limit applicability to clinical anxiety, other research indicates that findings from clinical  
349 samples replicate in samples recruited online<sup>6,8</sup>. Furthermore, it is increasingly recognised that  
350 clinical disorders lie on a continuum from health to disorder<sup>48</sup>. Although we did not deliberately  
351 set out to recruit individuals with clinically significant anxiety, 36% of our sample scored at or  
352 above a threshold designed for the detection of anxiety disorders on our measure of trait  
353 anxiety (see supplementary material). In light of this, and given limitations with research in  
354 clinical samples that includes medication load<sup>49</sup> and recruitment challenges<sup>50</sup>, online samples  
355 provide an effective method for studying clinically-relevant phenomena. Additionally, it is  
356 important to note that the effects we observed were small, as in previous studies using large-  
357 scale online testing<sup>6,25,25</sup>. However, large samples provide accurate effect size estimates in  
358 contrast to the exaggerated effects that are common in studies using small samples<sup>51</sup>. Such  
359 small effects are unsurprising given the multifactorial nature of psychiatric disorders<sup>52</sup>. While  
360 we have shown aversive learning to be important, we acknowledge this is likely to be one of a  
361 multitude of processes involved in the development of these conditions.

362 In conclusion, our results demonstrate links between transdiagnostic symptoms of psychiatric  
363 disorders and mechanisms of threat learning and uncertainty estimation in aversive  
364 environments. The findings emphasise the importance of these processes not only in anxiety  
365 but indicate a likely relevance across a spectrum of psychopathology.

## 366 Methods

### 367 Ethics

368 This research was approved by the University College London research ethics committee  
369 (reference 9929/003). All participants provided informed consent and were compensated  
370 financially for their time at a rate of at least £6 per hour.

### 371 Participants

372 We recruited 400 participants through Prolific<sup>31</sup>. Subjects were selected based on being aged  
373 18-65 and having at least a 90% approval rate across studies they had previously participated  
374 in. As described in our preregistration, we used a precision-based stopping rule to determine  
375 our sample size, stopping at the point at which either the 95% highest posterior density interval  
376 (HPDI) for all effects in our regression model reached 0.15 (checking with each 50 subjects  
377 recruited) or we had recruited 400 subjects. The precision target was not reached, and so we  
378 stopped at 400 subjects.

## Aversive learning and psychopathology

### 379 Avoidance learning task

380 Traditional lab-based threat learning tasks typically use aversive stimuli such as electric shocks  
381 as outcomes to be avoided. As it is not possible to use these stimuli online, we developed a  
382 game-based task in which subjects' goal was to avoid negative outcomes. While no primary  
383 aversive stimuli were used, and subjects received no actual monetary reward, there is an  
384 extensive literature showing that video games without such outcomes evoke strong positive  
385 and negative emotional experiences<sup>43-47</sup>, making this a promising method for designing an  
386 aversive learning task. In this game, participants were tasked with flying a spaceship through  
387 asteroid belts. Subjects were able to move the spaceship in the Y-axis alone, and this resulted  
388 in a one dimensional behavioural output. Crashing into asteroids diminished the spaceship's  
389 integrity by 10%. The spaceship's integrity slowly increased over the course of the task,  
390 however if enough asteroids were hit the integrity reduced to zero and the game finished. In  
391 this eventuality subjects were able to restart and continue where they left off. The overarching  
392 goal was to maximise the number of points scored, where the latter accumulated continuously  
393 for as long as the game was ongoing, and reset if the spaceship was destroyed. Subjects were  
394 shown the current integrity of the spaceship by a bar displayed in the corner of the screen,  
395 along with by a display of their current score.

396 Crucially, the location of safe spaces in the asteroid belts could be learned, and learning  
397 facilitated performance as it allowed correct positioning of the spaceship prior to observing the  
398 safe location. The task was designed such that without such pre-emptive positioning it was  
399 near impossible to successfully avoid the asteroids, thus encouraging subjects to learn the  
400 safest positions. Holes in the asteroids could appear either at the top or bottom of the screen  
401 (Figure 1A), and the probability of safety associated with either location varied independently  
402 over the course of the task. Thus, it was possible to learn the safety probability associated with  
403 each safety zone and adapt one's behaviour accordingly. The probability of each zone being  
404 safe was largely independent from the other (so that observing safety in one zone did not  
405 necessarily indicate the other was dangerous), although at least one zone was always safe on  
406 each trial. Participants also completed a control task that required avoidance that was not  
407 dependent on learning, enabling us to control for general motor-related avoidance ability in  
408 further analyses (described in supplementary material). We elected a priori to exclude subjects  
409 with limited response variability (indicated by a standard deviation of their positions below  
410 0.05) so as to remove subjects who did not move the spaceship. However, no subject met this  
411 exclusion criterion.

412 After completing the task, subjects were asked to provide ratings indicating how anxious the  
413 task made them feel and how motivated they were to avoid the asteroids, using visual analogue  
414 scales ranging from 0 to 100.

### 415 Behavioural data extraction

416 For analysis, we treated each pass through an asteroid belt as a trial. Overall there were 269  
417 trials in total. As a measure of behaviour, we extracted the mean Y position across the 1 second  
418 prior to observing the asteroid belt, representing where subjects were positioning themselves  
419 in preparation for the upcoming asteroid belt. This Y position was used for subsequent model



## Aversive learning and psychopathology

420 fitting. On each trial, the outcome for each zone was regarded as “danger” if asteroids were  
421 observed (regardless of whether they were hit by the subject) or safety if a hole in the asteroid  
422 belt was observed.

### 423 Computational modelling of behaviour

424 Our modelling approach focused on models that allowed the quantification of subjective  
425 uncertainty. To this end, we modelled behaviour using approximate Bayesian models that  
426 assume subjects estimate safety probability using a beta distribution. This approach is naturally  
427 suited to probability estimation tasks, as the beta distribution is bounded between zero and  
428 one, and provides a measure of uncertainty through the variance of the distribution. While  
429 certain reinforcement learning formulations can achieve similar uncertainty-dependent learning  
430 and quantification of uncertainty, we chose beta models as they have an advantage of being  
431 computationally simple. Empirically, these models have been used successfully in previous  
432 studies to capture value-based learning<sup>53</sup>, where they explain behaviour in aversive learning  
433 tasks better than commonly used reinforcement learning models<sup>15,54</sup>, a pertinent characteristic  
434 in the current task.

435 The basic premise underlying these models is that evidence for a given outcome is dependent  
436 on the number of times this outcome has occurred previously. For example, evidence for safety  
437 in a given location should then be highest when safety has been encountered many times in  
438 this location. This count can be represented by a parameter  $A$ , which is then incremented by a  
439 given amount every time safety is encountered. Danger is represented by a complementary  
440 parameter  $B$ . The balance between these parameters provides an indication of which outcome  
441 is most likely. Meanwhile, the overall number of outcomes counted influences the variance of  
442 the distribution and hence the uncertainty about this estimate. Thus, uncertainty is highest  
443 when few outcomes have been observed. The exact amount by which  $A$  and  $B$  are updated  
444 after every observed outcome can be estimated as a free parameter (here termed  $\tau$ ), and we  
445 can build asymmetry in learning into the model, so that learning about safety and danger have  
446 different rates, allowing updates for  $A$  and  $B$  to take on different values (here termed  $\tau^+$  and  $\tau^-$   
447 ).

448 Such a model is appropriate in stationary environments, when the probability of a given  
449 outcome is assumed to be constant throughout the experiment. However, in our task the  
450 probability of safety varied, and so it was necessary to build a forgetting process into the model.  
451 This is achieved by incorporating a decay (represented by parameter  $\lambda$ ) which diminishes the  
452 current values of  $A$  and  $B$  on every trial. The result of this process is akin to reducing the number  
453 of times they have been observed, and maintains the model’s ability to update in response to  
454 incoming evidence. It would also be possible to build asymmetry into the model here, where  
455 subjects could forget about positive and negative outcomes at different rate. However, testing  
456 this model in pilot data revealed that separate decay rates for each valence were not  
457 recoverable. Estimates for  $A$  and  $B$  are therefore updated on each trial ( $t$ ) according to the  
458 following equation for both safety zones independently (termed  $X$  and  $Y$  here). Both zones are  
459 updated on every trial, as subjects saw the outcome associated with both simultaneously. This  
460 formed the basis of all the probabilistic models tested:

## Aversive learning and psychopathology

$$461 \quad A_{t+1}^X = (1 - \lambda) \cdot A_t^X + \text{outcome}_t^X \cdot \tau^+ \cdot W \quad (1)$$

$$462 \quad B_{t+1}^X = (1 - \lambda) \cdot B_t^X + (1 - \text{outcome}_t^X) \cdot \tau^- \cdot W \quad (2)$$

463 We also observed in pilot data that subjects tended to be influenced more by outcomes  
464 occurring in the zone they had previously chosen, an effect likely due to attention. On this basis,  
465 we incorporated a weighting parameter that allowed the outcome of the unchosen option to  
466 be down-weighted by an amount shown in the above equation ( $W$ ) determined by an  
467 additional free parameter,  $\omega$ .

$$468 \quad W_{t+1}^X = \begin{cases} 1 & \text{if chosen} \\ \omega & \text{if unchosen} \end{cases} \quad (3)$$

469 We can calculate the estimated safety probability for each zone ( $P$ ) by taking the mean of this  
470 distribution:

$$471 \quad P_{t+1}^X = \frac{A_{t+1}^X}{(A_{t+1}^X + B_{t+1}^X)} \quad (4)$$

472 Similarly, we can derive a measure of uncertainty on each trial by taking the variance of this  
473 distribution.

$$474 \quad \sigma_{t+1}^X = \frac{A_{t+1}^X \cdot B_{t+1}^X}{(A_{t+1}^X + B_{t+1}^X)^2 \cdot (A_{t+1}^X + B_{t+1}^X + 1)} \quad (5)$$

475

476 In order to fit our model to the observed behaviour, we require an output that represents the  
477 position of the spaceship on the screen. This position ( $pos$ ) was calculated based on the safety  
478 probability of the two safety zones, such that the position was biased towards the safest  
479 location and was nearer the centre of the screen when it was unclear which position was safest.

$$480 \quad pos_{t+1} = \frac{(P_{t+1}^X - P_{t+1}^Y) + 1}{2} \quad (6)$$

481 Further models elaborated on this basic premise, and full details are provided in supplementary  
482 material. For completeness, we also tested two reinforcement learning models, a Rescorla-  
483 Wagner model and a variant of this model with different learning rates for better and worse  
484 than expected outcomes<sup>32</sup>, both of which are described in supplementary material. However,  
485 we focus on the probabilistic models due to their ability to represent uncertainty naturally; our  
486 primary aim was not to differentiate between probabilistic and reinforcement learning models,  
487 but to use previously validated models to provide insights into the relationship between  
488 aversive learning, uncertainty, and psychopathology.

489 Models were fit with a hierarchical Bayesian approach using variational inference implemented  
490 in PyMC3, through maximising the likelihood of the data given a reparametrised beta  
491 distribution with a mean provided by the model and a single free variance parameter. Model fit  
492 was assessed using the Watanabe-Akaike Information Criterion (WAIC)<sup>55</sup>, an index of model

## Aversive learning and psychopathology

493 fit designed for Bayesian models that accounts for model complexity. Parameter distributions  
494 were visualised using raincloud plots<sup>56</sup>.

### 495 Measures of psychiatric symptoms

496 Our first set of hypotheses focused on state/trait anxiety and intolerance of uncertainty. These  
497 were measured using the State Trait Inventory of Cognitive and Somatic Anxiety (STICSA)<sup>57</sup>  
498 and the Intolerance of Uncertainty Scale (IUS)<sup>58</sup> respectively. We also wished to examine how  
499 behaviour in our task related to the three transdiagnostic factors identified by Gillan et al.  
500 (2016), based on factor analysis of a range of psychiatric measures. To measure these factors  
501 more efficiently, we developed a reduced set of questions that provided an accurate  
502 approximation of the true factor scores, details of which are provided in supplementary  
503 material.

### 504 Regression models

505 Bayesian regression models were used to investigate relationships between behaviour and  
506 psychiatric measures, predicting each behavioural measure of interest from the psychiatric  
507 measures. Our dependent variables were parameters and quantities derived from our model,  
508 which represented the way in which an individual learns about safety probability and how they  
509 estimate uncertainty. Specifically, we used the two update parameters from our model ( $\tau^+$  and  
510  $\tau^-$ , referring to the extent to which subjects update in response to safety and danger  
511 respectively) and the mean safety probability and uncertainty estimates across the task  
512 (generated by simulating data from the model with each subject's estimated parameter values).  
513 Crucially, the fact that task outcomes were identical for every subject ensured these values  
514 were dependent only on the manner by which subjects learned about safety, not the task itself.

515 These models were constructed using Bambi<sup>59</sup> and fit using Markov chain Monte Carlo (MCMC)  
516 sampling, each with 8000 samples, 2000 of which were used for burn-in. All models included  
517 age and sex as covariates, along with performance on our control task to account for non-  
518 learning related avoidance ability. For analyses predicting state and trait anxiety and  
519 intolerance of uncertainty, we constructed a separate model for each variable due to the high  
520 collinearity between these measures. For analyses including the three transdiagnostic factors,  
521 these were entered into a single model. When reporting regression coefficients, we report the  
522 mean of the posterior distribution along with the 95% highest posterior density interval (HPDI),  
523 representing the points between which 95% of the posterior distribution's density lies. All  
524 analyses were specified in our preregistration. We did not correct for multiple comparisons in  
525 these analyses as our approach uses Bayesian parameter estimation, rather than frequentist  
526 null hypothesis significance testing, and as such multiple comparison correction is unnecessary  
527 and incompatible with this method<sup>60</sup>.

528

### 529 Partial least squares regression

530 To provide a data-driven characterisation of the relationship between task behaviour and  
531 psychiatric symptoms, and identify transdiagnostic components that are grounded in both self-  
532 report and behaviour, we used partial least squares (PLS) regression to identify dimensions of

## Aversive learning and psychopathology

533 covariance between individual questions and the measures derived from our modelling. We  
534 excluded the STICSA state subscale from this analysis, so that only trait measures were  
535 included. To ensure robustness of these results, we split our data into training and testing sets,  
536 made up of 75% and 25% of the data respectively. To identify the appropriate number of  
537 components within the training set, we used a 10-fold cross-validation procedure, fitting the  
538 model on 90% of the training data and evaluating its performance on the left-out 10%. The  
539 mean squared error of the model's predictions was then averaged across test folds to provide  
540 an index of the model's predictive accuracy with different numbers of components, using cross-  
541 validation to reduce the risk of overfitting

542 Once the number of components was determined, we validated the model's predictions by  
543 testing its predictive accuracy on the held-out 25% of the data. To provide a measure of  
544 statistical significance we used permutation testing, fitting the model on the training data 1000  
545 times with shuffled outcome variables and then testing each fitted model on the held-out data,  
546 to assess its predictive accuracy when fitted on data where no relationship exists between the  
547 predictors and outcomes. This procedure provides a null distribution, from which we can then  
548 determine the likelihood of observing predictive accuracy at least as high as that found in the  
549 true data under the null hypothesis.

550 Recent work has highlighted the risks inherent in PLS-like methods when used in high  
551 dimensional datasets<sup>36</sup>, namely that they can easily be overfit resulting in solutions that do not  
552 generalise beyond the data used to fit the model. Our approach avoids these problems by  
553 evaluating the performance on our model 25% of the data that has been held out from the  
554 model fitting stage.

## 555 Preregistration and data availability

556 The main hypotheses and methods of this study were preregistered on the Open Science  
557 Framework (<https://osf.io/jp5qn>). The data-driven PLS regression analysis was exploratory.  
558 Data is available at <https://osf.io/b95w2/> and code is available at  
559 <https://github.com/tobywise/online-aversive-learning>.

560

## 561 Acknowledgements

562 T.W. is supported by a Wellcome Trust Sir Henry Wellcome Fellowship (206460/17/Z). R.J.D.  
563 holds a Wellcome Trust Investigator award (098362/Z/12/Z). The Max Planck UCL Centre is a  
564 joint initiative supported by UCL and the Max Planck Society. The Wellcome Centre for Human  
565 Neuroimaging is supported by core funding from the Wellcome Trust (203147/Z/16/Z). We  
566 thank Evan Russek for comments on an earlier version of this manuscript.

## 567 Disclosures

568 The authors report no conflicts of interest in relation to this work.

569

## Aversive learning and psychopathology

### 570 References

- 571 1. Reininghaus, U. et al. Stress Sensitivity, Aberrant Salience, and Threat Anticipation in  
572 Early Psychosis: An Experience Sampling Study. *Schizophr Bull* **42**, 712–722 (2016).
- 573 2. Harrison, A., Tchanturia, K. & Treasure, J. Attentional Bias, Emotion Recognition, and  
574 Emotion Regulation in Anorexia: State or Trait? *Biological Psychiatry* **68**, 755–761 (2010).
- 575 3. Friston, K. J., Stephan, K. E., Montague, R. & Dolan, R. J. Computational psychiatry: the  
576 brain as a phantastic organ. *The Lancet Psychiatry* **1**, 148–158 (2014).
- 577 4. Montague, P. R., Dolan, R. J., Friston, K. J. & Dayan, P. Computational psychiatry. *Trends in*  
578 *Cognitive Sciences* **16**, 72–80 (2012).
- 579 5. Browning, M., Behrens, T. E., Jocham, G., O'Reilly, J. X. & Bishop, S. J. Anxious individuals  
580 have difficulty learning the causal statistics of aversive environments. *Nat Neurosci* **18**,  
581 590–596 (2015).
- 582 6. Gillan, C. M., Kosinski, M., Whelan, R., Phelps, E. A. & Daw, N. D. Characterizing a  
583 psychiatric symptom dimension related to deficits in goal-directed control. *eLife* **5**,  
584 e11305 (2016).
- 585 7. Huang, H., Thompson, W. & Paulus, M. P. Computational Dysfunctions in Anxiety:  
586 Failure to Differentiate Signal From Noise. *Biological Psychiatry* **82**, 440–446 (2017).
- 587 8. Seow, T. X. & Gillan, C. Transdiagnostic phenotyping reveals a range of metacognitive  
588 deficits associated with compulsivity. *bioRxiv* 664003 (2019) doi:10.1101/664003.
- 589 9. Lawson, R. P., Mathys, C. & Rees, G. Adults with autism overestimate the volatility of the  
590 sensory environment. *Nat Neurosci* **advance online publication**, (2017).

## Aversive learning and psychopathology

- 591 10. Mkrtchian, A., Aylward, J., Dayan, P., Roiser, J. P. & Robinson, O. J. Modeling Avoidance in  
592 Mood and Anxiety Disorders Using Reinforcement Learning. *Biological Psychiatry* **82**,  
593 532–539 (2017).
- 594 11. Vinckier, F. et al. Confidence and psychosis: a neuro-computational account of  
595 contingency learning disruption by NMDA blockade. *Mol Psychiatry* **21**, 946–955 (2016).
- 596 12. Aylward, J. et al. Altered learning under uncertainty in unmedicated mood and anxiety  
597 disorders. *Nature Human Behaviour* **1** (2019) doi:10.1038/s41562-019-0628-0.
- 598 13. Robinson, O. J. & Chase, H. W. Learning and Choice in Mood Disorders: Searching for the  
599 Computational Parameters of Anhedonia. *Computational Psychiatry* 1–26 (2017)  
600 doi:10.1162/CPSY\_a\_00009.
- 601 14. Sharp, P. B. & Eldar, E. Computational Models of Anxiety: Nascent Efforts and Future  
602 Directions. *Curr Dir Psychol Sci* **28**, 170–176 (2019).
- 603 15. Wise, T., Michely, J., Dayan, P. & Dolan, R. J. A computational account of threat-related  
604 attentional bias. *PLOS Computational Biology* **15**, e1007341 (2019).
- 605 16. Pulcu, E. & Browning, M. The Misestimation of Uncertainty in Affective Disorders. *Trends*  
606 *in Cognitive Sciences* (2019) doi:10.1016/j.tics.2019.07.007.
- 607 17. Bach, D. R. & Dolan, R. J. Knowing how much you don't know: a neural organization of  
608 uncertainty estimates. *Nat Rev Neurosci* **13**, 572–586 (2012).
- 609 18. de Berker, A. O. et al. Computations of uncertainty mediate acute stress responses in  
610 humans. *Nat Commun* **7**, 10996 (2016).
- 611 19. Mathys, C., Daunizeau, J., Friston, K. J. & Stephan, K. E. A Bayesian foundation for  
612 individual learning under uncertainty. *Front. Hum. Neurosci.* **5**, 39 (2011).

## Aversive learning and psychopathology

- 613 20. Mathys, C. et al. Uncertainty in perception and the Hierarchical Gaussian Filter. *Front.*  
614 *Hum. Neurosci* **8**, 825 (2014).
- 615 21. Grupe, D. W. & Nitschke, J. B. Uncertainty and anticipation in anxiety: an integrated  
616 neurobiological and psychological perspective. *Nat Rev Neurosci* **14**, 488–501 (2013).
- 617 22. Dugas, M. J., Gagnon, F., Ladouceur, R. & Freeston, M. H. Generalized anxiety disorder: a  
618 preliminary test of a conceptual model. *Behaviour Research and Therapy* **36**, 215–226  
619 (1998).
- 620 23. Barlow, D. H. Unraveling the mysteries of anxiety and its disorders from the perspective  
621 of emotion theory. *Am Psychol* **55**, 1247–1263 (2000).
- 622 24. Grillon, C. Associative learning deficits increase symptoms of anxiety in humans.  
623 *Biological Psychiatry* **51**, 851–858 (2002).
- 624 25. Patzelt, E. H., Kool, W., Millner, A. J. & Gershman, S. J. Incentives Boost Model-Based  
625 Control Across a Range of Severity on Several Psychiatric Constructs. *Biological*  
626 *Psychiatry* (2018) doi:10.1016/j.biopsych.2018.06.018.
- 627 26. Rouault, M., Seow, T., Gillan, C. M. & Fleming, S. M. Psychiatric Symptom Dimensions Are  
628 Associated With Dissociable Shifts in Metacognition but Not Task Performance.  
629 *Biological Psychiatry* **84**, 443–451 (2018).
- 630 27. Norbury, A., Robbins, T. W. & Seymour, B. Value generalization in human avoidance  
631 learning. *eLife* **7**, e34779 (2018).
- 632 28. Michelini, G. et al. Delineating and validating higher-order dimensions of psychopathology  
633 in the Adolescent Brain Cognitive Development (ABCD) study. *Transl Psychiatry* **9**, 1–15  
634 (2019).

## Aversive learning and psychopathology

- 635 29. Blanco, C. et al. Mapping Common Psychiatric Disorders: Structure and Predictive Validity  
636 in the National Epidemiologic Survey on Alcohol and Related Conditions. *JAMA*  
637 *Psychiatry* **70**, 199–207 (2013).
- 638 30. Røyamb, E. et al. The joint structure of DSM-IV Axis I and Axis II disorders. *Journal of*  
639 *Abnormal Psychology* **120**, 198–209 (2011).
- 640 31. Palan, S. & Schitter, C. Prolific.ac—A subject pool for online experiments. *Journal of*  
641 *Behavioral and Experimental Finance* **17**, 22–27 (2018).
- 642 32. Lefebvre, G., Lebreton, M., Meyniel, F., Bourgeois-Gironde, S. & Palminteri, S. Behavioural  
643 and neural characterization of optimistic reinforcement learning. *Nature Human*  
644 *Behaviour* **1**, 0067 (2017).
- 645 33. Kebets, V. et al. Somatosensory-Motor Dysconnectivity Spans Multiple Transdiagnostic  
646 Dimensions of Psychopathology. *Biological Psychiatry* **86**, 779–791 (2019).
- 647 34. Jessen, K. et al. Patterns of Cortical Structures and Cognition in Antipsychotic-Naïve  
648 Patients With First-Episode Schizophrenia: A Partial Least Squares Correlation Analysis.  
649 *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* **4**, 444–453 (2019).
- 650 35. Mihalik, A. et al. Multiple hold-outs with stability: improving the generalizability of  
651 machine learning analyses of brain-behaviour relationships: A novel framework to link  
652 behaviour to neurobiology. *Biological Psychiatry* (2019)  
653 doi:10.1016/j.biopsych.2019.12.001.
- 654 36. Dinga, R. et al. Evaluating the evidence for biotypes of depression: Methodological  
655 replication and extension of Drysdale et al. (2017). *NeuroImage: Clinical* **22**, 101796  
656 (2019).



## Aversive learning and psychopathology

- 657 37. Plichta, M. M. & Scheres, A. Ventral–striatal responsiveness during reward anticipation in  
658 ADHD and its relation to trait impulsivity in the healthy population: A meta-analytic  
659 review of the fMRI literature. *Neuroscience & Biobehavioral Reviews* **38**, 125–134  
660 (2014).
- 661 38. Deakin, J., Aitken, M., Robbins, T. & Sahakian, B. J. Risk taking during decision-making in  
662 normal volunteers changes with age. *Journal of the International Neuropsychological*  
663 *Society* **10**, 590–598 (2004).
- 664 39. Pulcu, E. & Browning, M. Affective bias as a rational response to the statistics of rewards  
665 and punishments. *eLife Sciences* **6**, e27879 (2017).
- 666 40. Michely, J., Rigoli, F., Rutledge, R. B., Hauser, T. U. & Dolan, R. J. Distinct Processing of  
667 Aversive Experience in Amygdala Subregions. *Biological Psychiatry: Cognitive*  
668 *Neuroscience and Neuroimaging* (2019) doi:10.1016/j.bpsc.2019.07.008.
- 669 41. Beesdo-Baum, K. et al. Avoidance, Safety Behavior, and Reassurance Seeking in  
670 Generalized Anxiety Disorder. *Depression and Anxiety* **29**, 948–957 (2012).
- 671 42. Salkovskis, P. M. The Importance of Behaviour in the Maintenance of Anxiety and Panic: A  
672 Cognitive Account. *Behavioural and Cognitive Psychotherapy* **19**, 6–19 (1991).
- 673 43. Schneider, E. F., Lang, A., Shin, M. & Bradley, S. D. Death with a Story: How Story Impacts  
674 Emotional, Motivational, and Physiological Responses to First-Person Shooter Video  
675 Games. *Hum Commun Res* **30**, 361–375 (2004).
- 676 44. Ravaja, N., Saari, T., Salminen, M., Laarni, J. & Kallinen, K. Phasic Emotional Reactions to  
677 Video Game Events: A Psychophysiological Investigation. *Media Psychology* **8**, 343–367  
678 (2006).

## Aversive learning and psychopathology

- 679 45. Ravaja, N., Turpeinen, M., Saari, T., Puttonen, S. & Keltikangas-Järvinen, L. The  
680 psychophysiology of James Bond: Phasic emotional responses to violent video game  
681 events. *Emotion* **8**, 114–120 (2008).
- 682 46. Hutton, E. & Sundar, S. S. Can Video Games Enhance Creativity? Effects of Emotion  
683 Generated by Dance Dance Revolution. *Creativity Research Journal* **22**, 294–303 (2010).
- 684 47. Ravaja, N. et al. Spatial Presence and Emotions during Video Game Playing: Does It  
685 Matter with Whom You Play? *Presence: Teleoperators and Virtual Environments* **15**,  
686 381–392 (2006).
- 687 48. Plomin, R., Haworth, C. M. A. & Davis, O. S. P. Common disorders are quantitative traits.  
688 *Nat. Rev. Genet.* **10**, 872–878 (2009).
- 689 49. Phillips, M. L., Travis, M. J., Fagiolini, A. & Kupfer, D. J. Medication Effects in Neuroimaging  
690 Studies of Bipolar Disorder. *American Journal of Psychiatry* **165**, 313–320 (2008).
- 691 50. Wise, T. et al. Recruiting for research studies using online public advertisements:  
692 examples from research in affective disorders. *Neuropsychiatr Dis Treat* **12**, 279–285  
693 (2016).
- 694 51. Gelman, A. & Carlin, J. Beyond Power Calculations: Assessing Type S (Sign) and Type M  
695 (Magnitude) Errors. *Perspect Psychol Sci* **9**, 641–651 (2014).
- 696 52. Kendler, K. S. From Many to One to Many—the Search for Causes of Psychiatric Illness.  
697 *JAMA Psychiatry* **76**, 1085–1091 (2019).
- 698 53. de Boer, L. et al. Attenuation of dopamine-modulated prefrontal value signals underlies  
699 probabilistic reward learning deficits in old age. *eLife* **6**, e26424 (2017).
- 700 54. Tzovara, A., Korn, C. W. & Bach, D. R. Human Pavlovian fear conditioning conforms to  
701 probabilistic learning. *PLOS Computational Biology* **14**, e1006243 (2018).

## Aversive learning and psychopathology

- 702 55. Watanabe, S. Asymptotic Equivalence of Bayes Cross Validation and Widely Applicable  
703 Information Criterion in Singular Learning Theory. *J. Mach. Learn. Res.* **11**, 3571–3594  
704 (2010).
- 705 56. Allen, M., Poggiali, D., Whitaker, K., Marshall, T. R. & Kievit, R. A. Raincloud plots: a multi-  
706 platform tool for robust data visualization. *Wellcome Open Res* **4**, (2019).
- 707 57. Grös, D. F., Antony, M. M., Simms, L. J. & McCabe, R. E. Psychometric properties of the  
708 State-Trait Inventory for Cognitive and Somatic Anxiety (STICSA): comparison to the  
709 State-Trait Anxiety Inventory (STAI). *Psychol Assess* **19**, 369–381 (2007).
- 710 58. Carleton, R. N., Norton, M. A. P. J. & Asmundson, G. J. G. Fearing the unknown: A short  
711 version of the Intolerance of Uncertainty Scale. *Journal of Anxiety Disorders* **21**, 105–117  
712 (2007).
- 713 59. Yarkoni, T. & Westfall, J. Bambi: A simple interface for fitting Bayesian mixed effects  
714 models. (2016) doi:10.31219/osf.io/rv7sn.
- 715 60. Kruschke, J. K. & Liddell, T. M. The Bayesian New Statistics: Hypothesis testing,  
716 estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychon Bull*  
717 *Rev* **25**, 178–206 (2018).
- 718