

 Open access • Posted Content • DOI:10.1101/843045

## **Associations between aversive learning processes and transdiagnostic psychiatric symptoms revealed by large-scale phenotyping** — [Source link](#)

Toby Wise, Raymond J. Dolan

**Institutions:** University College London

**Published on:** 15 Nov 2019 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Social anxiety, Anxiety and Impulsivity

Related papers:

- [Associations between aversive learning processes and transdiagnostic psychiatric symptoms in a general population sample](#)
- [Transdiagnostic Phenotyping Reveals a Host of Metacognitive Deficits Implicated in Compulsivity](#)
- [Modeling Avoidance in Mood and Anxiety Disorders Using Reinforcement Learning](#)
- [Decision-Making in Anxiety and Its Disorders](#)
- [Modeling Trait Anxiety: From Computational Processes to Personality](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/associations-between-aversive-learning-processes-and-387qguzo8u>

# Associations between aversive learning processes and transdiagnostic psychiatric symptoms in a general population sample

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

Symptom expression in psychiatric conditions is often linked to altered threat perception, however how computational mechanisms that support aversive learning relate to specific psychiatric symptoms remains undetermined. We answer this question using an online game-based aversive learning task together with measures of common psychiatric symptoms in 400 subjects. We show that physiological symptoms of anxiety and a transdiagnostic compulsivity-related factor are associated with enhanced safety learning, as measured using a probabilistic computational model, while trait cognitive anxiety symptoms are associated with enhanced learning from danger. We use data-driven partial least squares regression to identify two separable components across behavioural and questionnaire data: one linking enhanced safety learning and lower estimated uncertainty to physiological anxiety, compulsivity, and impulsivity; the other linking enhanced threat learning and heightened uncertainty estimation to symptoms of depression and social anxiety. Our findings implicate aversive learning processes in the expression of psychiatric symptoms that transcend diagnostic boundaries.

<sup>1</sup> Wellcome Centre for Human Neuroimaging, University College London, London, UK. <sup>2</sup> Max Planck UCL Centre for Computational Psychiatry and Ageing Research, University College London, London, UK. <sup>3</sup> Division of the Humanities and Social Sciences, California Institute of Technology, Pasadena, CA, USA.  
✉email: [t.wise@ucl.ac.uk](mailto:t.wise@ucl.ac.uk)

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

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

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

Existing work on aversive learning has had a particular focus on symptoms of anxiety and depression<sup>7,12</sup>. However, these approaches have not been designed optimally for identifying mechanisms that span traditional diagnostic boundaries. This assumes importance in light of recent studies, using large samples, showing several aspects of learning and decision-making

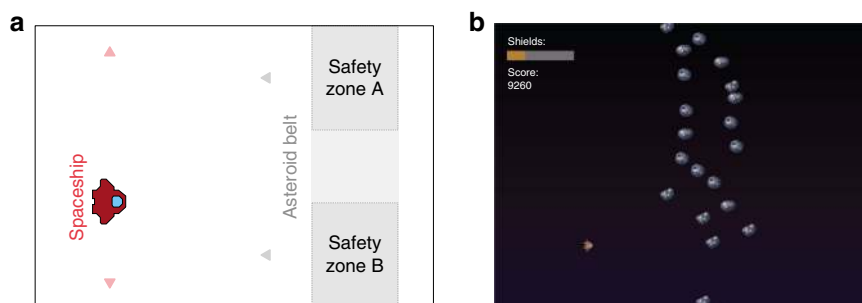
relate more strongly to transdiagnostic factors (symptom dimensions that are not unique to any one disorder) than to any specific categorical conception of psychiatric disorder<sup>6,8,25–27</sup>. Applying such an approach to aversive learning may yield better insights into the role of learning in psychiatric disorders. In addition, computationally defined measures of learning and decision-making can facilitate identification of novel transdiagnostic factors, going beyond those identified based solely on correlated symptom clusters in self-report and clinical interview measures<sup>6,28–30</sup>.

Here, we aim to clarify the nature of the relationship between aversive learning processes and traditional measures of anxiety, as well as transdiagnostic psychiatric factors identified in prior work<sup>6</sup> in a large, preregistered study conducted online. This allows us to measure effects with high precision, potentially helping to resolve mixed findings from previous studies<sup>12,15</sup>, in addition to identifying small but meaningful effects that cross traditional diagnostic boundaries<sup>6</sup>. As in similar prior studies<sup>6,8,25,26,31</sup>, we do not purposefully recruit subjects diagnosed with mental health problems, instead focusing on exploring relationships with the variation in symptoms present in the general population. While this does not allow concrete conclusions about clinical disorders per se, we note findings using similar approaches have replicated those seen in clinical samples<sup>8,32</sup>, providing reassurance these methods can generate insights into psychiatrically relevant phenomena. Thus, we use a computational approach to test whether anxiety and transdiagnostic symptoms are associated with biased learning from safety and threat, whether these factors relate to altered estimates of threat likelihood, and whether they are associated with different levels of uncertainty during threat learning. We then use partial least squares (PLS) regression, a data-driven multivariate method, to derive transdiagnostic latent components of psychopathology grounded in both self-report and computational measures. In contrast to our primary analyses, this analysis is exploratory and data-driven, enabling us to generate hypotheses for future research. Given difficulties in using traditional aversive stimuli in an online setting, we develop a game-based avoidance task designed to engage threat and avoidance processes without the need for administration of painful or noxious stimuli. Both the task and modelling are, in principle, similar to our previous lab-based task<sup>15</sup>, but their implementation here allows straightforward administration in large samples recruited online. Our results demonstrate that learning from safety and danger are associated with distinct symptom dimensions that cut across diagnostic boundaries, implicating aversive learning processes in a range of psychiatric symptoms.

## Results

**Task performance.** Four hundred subjects recruited online through Prolific<sup>33</sup> performed a game-based aversive learning task, where the aim was to fly a spaceship through asteroid belts without being hit (Fig. 1). Getting hit by the asteroids reduced the integrity of the spaceship, and after sufficient hits the game terminated. Crucially, there were two zones at the top and bottom of the screen where subjects could encounter a hole in the asteroid belt, each associated with a changing probability of being safe. In order to perform well at the task subjects needed to learn which zone was safest and behave accordingly.

Subjects were engaged and performed well at the task, with a median number of spaceship destructions of 1 (Interquartile range = 2) over the course of the task. They also reported high motivation to perform the task, providing a mean rating of 85.70 (SD = 18.44) when asked to rate how motivated they were to avoid asteroids on a scale from 0–100. Reassuringly, no subjects



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

met our exclusion criteria designed to remove those not engaging in the task.

**Computational modelling of behaviour.** To quantitatively describe behaviour, we fit a series of computational models to subjects' position data during the task (see “Methods” and Supplementary Methods for a full description of tested models). The winning model was a probabilistic model incorporating different updates parameters for safety and danger, as well as a stickiness parameter representing a tendency for subjects to stick with their previous position. This model represents an extension of one we have previously used successfully as a lab-based aversive learning task<sup>15</sup>, and is described fully in the “Methods” section. Briefly, this winning model assumes that subjects in the task represent the safety probability of each zone using a beta distribution, which is updated on each trial based on encounters with danger or safety. Simulating responses using the model, using each subject's estimated parameter values, produced behavioural profiles that demonstrated a high concordance with the true data, reproducing broad behavioural patterns seen in the true data (Supplementary Fig. 1). We note that we do not wish to make strong claims regarding the strategy employed by subjects in performing this task; our central focus was on exploring relationships with symptoms, and we focused on probabilistic models as they provide a natural measure of uncertainty that can be associated with symptoms. It is possible that more complex reinforcement learning models for example could produce equally good fits to the data.

**Task and model validation.** It was important to first ensure that the task has content validity, and that it produces behaviour reminiscent of more traditional tasks. Likewise, the computational models used should provide measures and parameter estimates that reflect the behaviour they aim to describe. We therefore conducted extensive validation exercises. These are reported fully in Supplementary Methods, but we summarise them here.

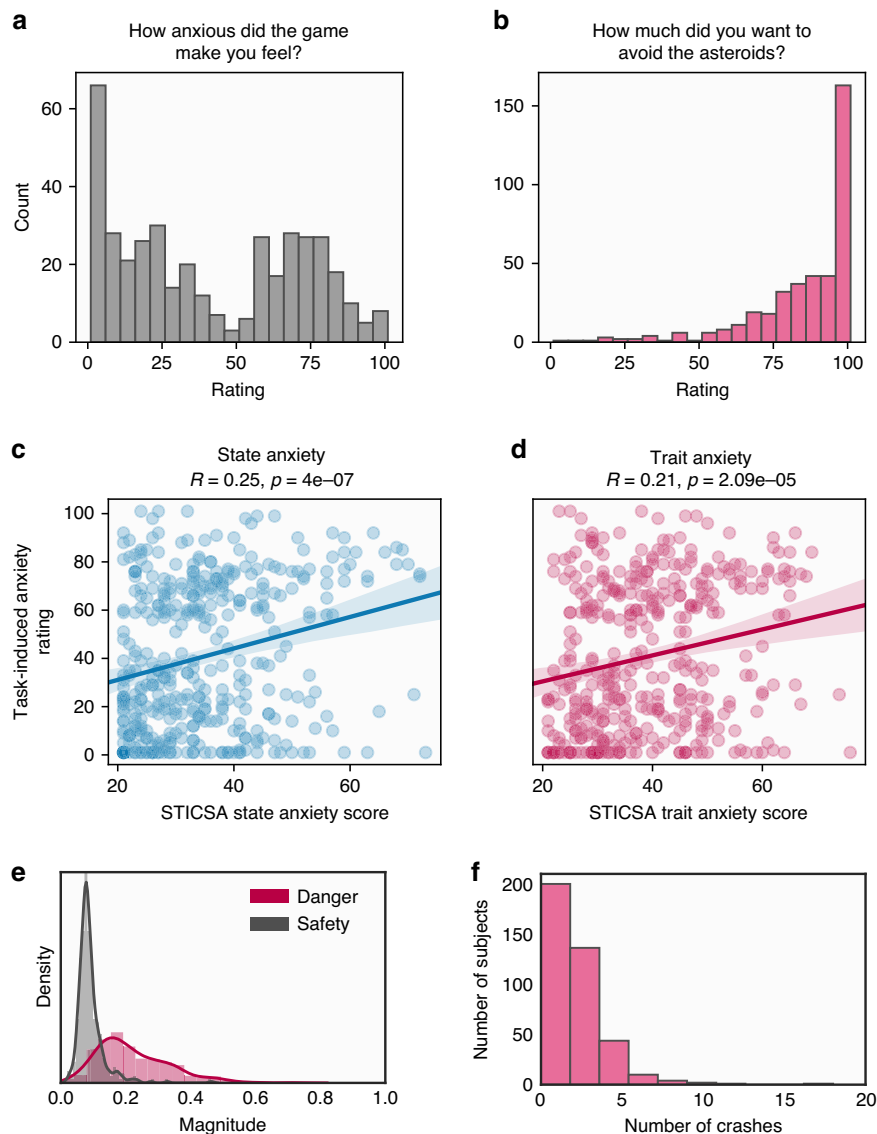
First, we ensured the task did induce states of subjective anxiety in the majority of subjects (Fig. 2a), and this level of anxiety was correlated with self-report state and trait anxiety (Fig. 2c, d). Importantly, subjects adjusted their position to a greater extent following danger than following safety (Fig. 2e), indicating that they were adapting their behaviour in response to outcomes in the task, rather than behaving randomly. With respect to our computational model, we verified that the model's update parameters were robustly correlated with subjects' tendency to move, or stay, following danger and safety, respectively (Fig. 3d).

We also assessed whether safety and uncertainty values, produced by simulating data from our model with best fitting parameters, related to subjects' model-free behaviour. We found that subjects changed their position more when model-derived uncertainty was high, and when the difference between the safety value of the two zones was small. This pattern (Fig. 3c) indicates that value and uncertainty measures do reflect meaningful quantities for behaviour. Finally, we verified that our model's update parameters showed greater updating from danger relative to safety, as we found in a previous lab-based study<sup>15</sup>, finding this was indeed the case (Fig. 3e).

Parameter recovery analyses are also reported in Supplementary Fig. 2 and indicated good recoverability (lowest  $r$  between true and recovered parameter value = 0.61). It should be noted that we did find a moderate negative correlation between the two update parameters ( $r = -0.49$ ), which may be an artefact of the fact that safety and danger outcomes in the task were partially anticorrelated (i.e., when one zone was likely to be safe, the other was likely to be dangerous). However, despite this, parameter recovery tests indicated good recoverability of these parameters.

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

However, it is possible that our state and trait anxiety sum scores may obscure more nuanced effects relating to different symptom dimensions. To test this, we performed exploratory analyses on the two subscales of our trait anxiety measure, which represent cognitive symptoms, such as worry, and physical symptoms, which reflect aspects of physiological arousal. This analysis revealed a dissociation between cognitive and somatic trait anxiety, whereby cognitive symptoms were associated with heightened learning from danger, heightened uncertainty, reduced learning from safety, and lower safety probability, while somatic anxiety showed an opposite pattern (Table 1 and Fig. 5). The same was true for state anxiety to a lesser extent; while effects for somatic anxiety remained strong, those for cognitive anxiety were weaker and had 95% HDPIs that



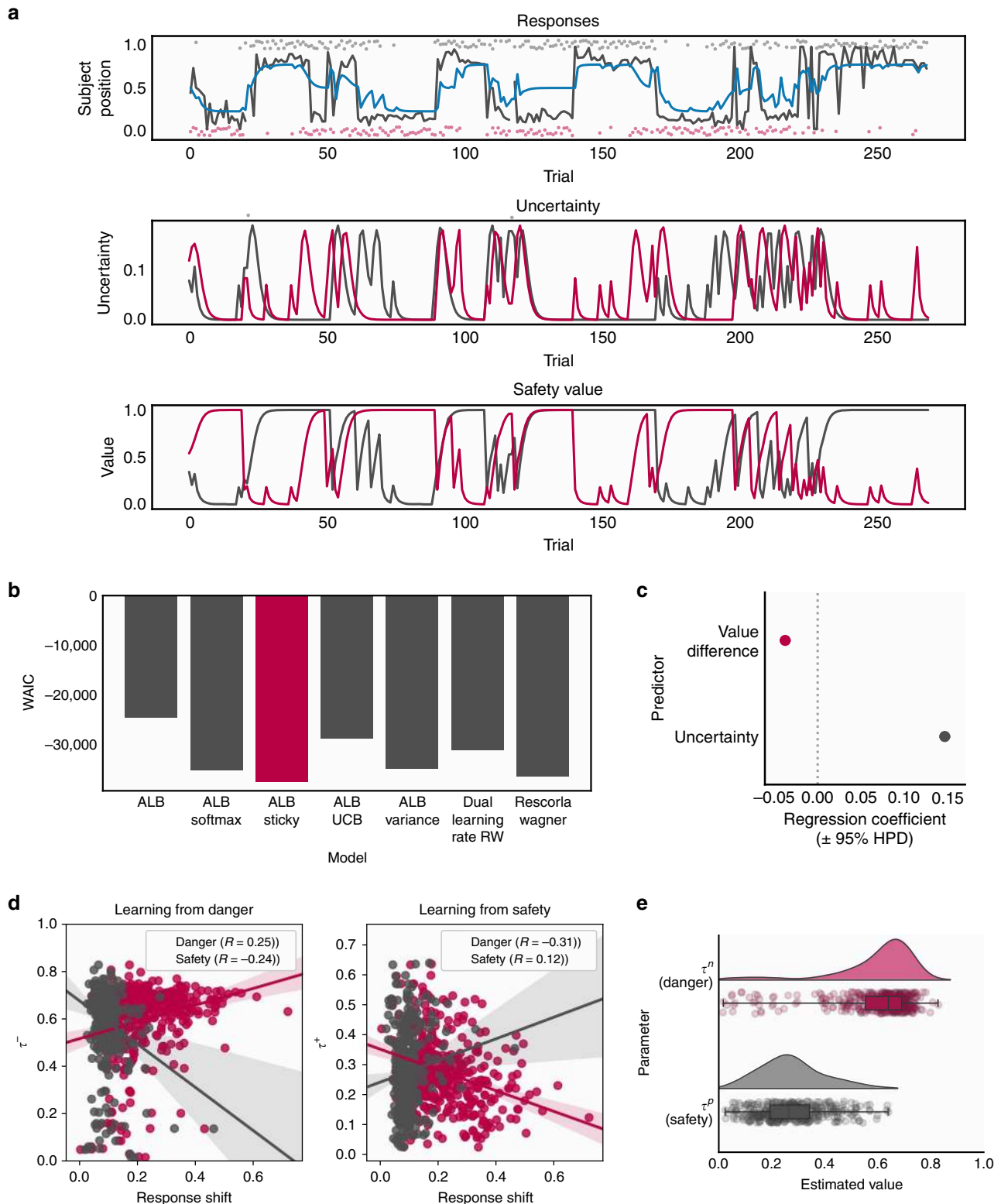
**Fig. 2 Subjective and behavioural responses during the task.** **a** Distribution of task-induced anxiety ratings recorded after the task. **b** Distribution of task motivation ratings. **c, d** Relationships between task-induced anxiety ratings and state and trait anxiety scores. Correlations represent Pearson  $r$  statistics, and  $p$  values are two sided without correction for multiple comparisons. **e** Degree of location switching after encountering danger and safety across subjects. The switch magnitude is the average absolute change in position between trial  $n$  and trial  $n + 1$ . As expected, subjects showed more switching behaviour after encountering danger and were more likely to stay in the same position following a safe outcome. **f** Distribution of crash number (representing the number of subjects hitting enough asteroids to end the game) across subjects.

included zero (Table 1 and Fig. 5). Together, these results indicate that different components of trait anxiety have distinct relationships with aversive learning processes.

**Associations between aversive learning and transdiagnostic factors.** Following this, we examined the extent to which task behaviour was associated with three transdiagnostic factors of psychopathology identified through self-report assessments in previous research<sup>6</sup>. Here, we observed effects of a factor labelled compulsivity and intrusive thought (Fig. 4, Table 2), reflecting the fact that subjects scoring higher on this factor learned faster about safety and had higher safety probability estimates. There was also a weak effect of this factor on uncertainty, although the HPDI for this included zero. Other effects were weak, and including reported task motivation as a covariate had a negligible effect on the results (see Supplementary Figs. 5 and 6). Importantly, all of these analyses were determined a priori and are included in our

preregistration. We also examined effects of age and sex, and their interactions with our primary variables. These results are described in Supplementary Fig. 4.

**Psychiatric constructs derived from behaviour and self-report.** Numerous studies have used dimensionality reduction procedures such as factor analysis on questionnaire-based data to identify factors of psychopathology that cut across diagnostic boundaries<sup>6,28–30</sup>. This, in turn, has revealed that many behaviourally defined phenotypes are more strongly associated with transdiagnostic factors than any single disorder<sup>6,8,26</sup>. We built upon this work by incorporating computationally derived indexes of behaviour into this dimensionality reduction procedure, where the aim was to identify latent constructs grounded in both self-report and behaviour. We used PLS regression, a method that identifies latent components linking multivariate data from multiple domains based on their shared covariance. This general



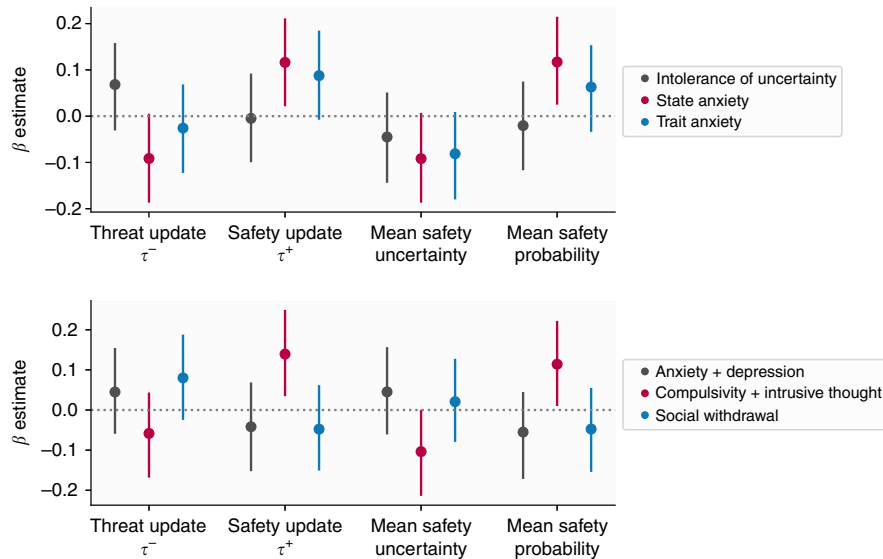
method has been employed successfully in prior studies to provide insight into how panels of cognitive and behavioural measures relate to multivariate neuroimaging-derived phenotypes<sup>34–36</sup>. PLS-like analyses can be problematic if not properly validated (for example producing spurious results due to overfitting), and so we adopted best-practice methods for validating these results<sup>37–39</sup>, selecting the optimal number of components using cross-

validation and training the model on 75% of the data, before testing its performance on the remaining 25% of the data. Importantly, in contrast to other analyses, this was an exploratory analysis where our aim was to provide indicative results that can be validated in future research.

We first identified the number of components that best describe our data by evaluating the performance of a predictive



**Fig. 3 Computational model fitting results.** **a** Data generated from the model. The top panel shows responses and model fit for an example subject. In the top panel, the grey line represents the subject's position throughout the task, with the grey and red dots representing safe locations on each trial. The blue line represents simulated data from the model for this subject. The lower two panels show estimated uncertainty and safety probability for each stimulus (represented by the grey and red lines) across the duration of the task, generated by simulating data from the model. **b** Model comparison results showing the Watanabe-Akaike Information Criterion (WAIC) score for each model with the winning model highlighted. ALB asymmetric leaky beta, RW Rescorla-Wagner. **c** Results of our analysis validating the safety value and uncertainty measures, showing the extent to which each measure predicted subjects' tendency to switch position (described in Supplementary Methods, with error bars estimated from 400 subjects). **d** Correlations between estimated update parameters for danger (left) and safety (right) and our behavioural measure of position switching after these outcomes across subjects, demonstrating that parameters from our model reflect purely behavioural characteristics. **e** Distributions of estimated parameter values for  $\tau^+$  and  $\tau^-$ , representing update rates following danger and safety outcomes, respectively, showing a bias in updating whereby subjects update to a greater extent in response to danger than safety. In the box plots, the centre of the box represents the median, the bounds of the box represent the quartiles (25 and 75%) of the distribution, and the whiskers represent the minimum and maximum values of the data.



**Fig. 4 Results from analyses relating model-derived measures to psychiatric symptoms.** Top panel: Results of state/trait anxiety and intolerance of uncertainty models, showing relationships between these psychiatric variables and behavioural variables. Points indicate the mean of the posterior distribution for the regression coefficient parameter, while error bars represent the 95% highest posterior density interval, estimated from 400 subjects. The  $\beta$  estimate here refers to the regression coefficient for each predictor. Bottom panel: Results of three factor model, showing relationships between behaviour and factors labelled anxiety and depression, compulsivity and intrusive thought, and social withdrawal.

PLS model using cross-validation. We found two latent components gave the best predictive performance (Fig. 6a). We then evaluated the performance of this model on held out data using permutation testing, showing our model achieved a statistically significant level of predictive accuracy (permutation  $p = 0.025$ , Fig. 6b). This indicates that our combined self-report and behavioural data is best explained by a two-component structure linking these two domains. Importantly, the fact that this level of accuracy was found on unseen data ensures that our results do not result from overfitting the training data<sup>37</sup>.

To aid interpretation of these two components we examined how behavioural variables loaded on each. The first component had positive weights on update rates in response to safety and estimated safety likelihood, and negative weights on update rates in response to threat, decay, stickiness, and mean uncertainty estimates (Fig. 6c), while the reverse was true of the second component. Loadings on questionnaire items were varied, and labelling such components is invariably subjective. Nevertheless, the first component tended to load most strongly on items describing physical symptoms of anxiety, compulsive behaviour, and impulsivity. In contrast, the second latent component loaded primarily on items describing social anxiety and depressed mood. For illustrative purposes, items with the top 10% percent of differences in loadings between components are shown in Fig. 6d, with full details available in Supplementary Fig. 9.

## Discussion

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

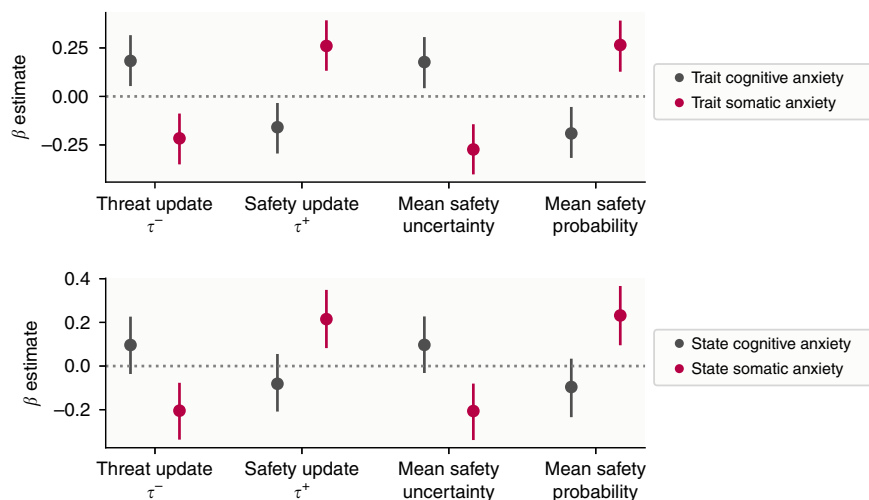
We found a counter-intuitive relationship between biases in learning and the presence of features of anxiety. Subjects scoring higher on state anxiety tended to update their predictions to a greater extent in response to safety, as well as perceiving safety to be more likely overall, than those scoring low on this measure. However, in exploratory analyses of anxiety subscales we observed a dissociation between cognitive and somatic symptoms of anxiety, whereby cognitive anxiety was associated with enhanced learning from threat, greater uncertainty, and lower safety probability estimates, while the reverse was true for somatic symptoms. Together, these results suggest divergent roles of aversive learning processes in distinct symptom dimensions.

Our results regarding state and trait anxiety as a whole diverge from previous findings that report individuals diagnosed with clinical anxiety and depression learn faster from punishment<sup>12</sup>, but are in concordance with our previous work in a non-clinical sample using a more traditional lab-based aversive learning task<sup>15</sup>. The large sample size employed here allowed us to

**Table 1 Estimates from regression model predicting learning-related variables derived from our computational model from measures of intolerance of uncertainty, state anxiety, and trait anxiety.**

Target variable	Predictor	Estimate (±95% HPDI)
Total scores		
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)
Trait anxiety subscales		
Threat update ( $\tau^-$ )	STICSA S <i>cognitive</i>	0.1 (−0.04, 0.23)
	STICSA S <i>somatic</i>	<b>−0.2 (−0.34, −0.08)</b>
	STICSA T <i>cognitive</i>	<b>0.18 (0.05, 0.32)</b>
Safety update ( $\tau^+$ )	STICSA T <i>somatic</i>	<b>−0.22 (−0.35, −0.09)</b>
	STICSA S <i>cognitive</i>	−0.08 (−0.21, 0.06)
	STICSA S <i>somatic</i>	<b>0.22 (0.08, 0.35)</b>
Mean safety uncertainty	STICSA T <i>cognitive</i>	<b>−0.16 (−0.29, −0.03)</b>
	STICSA T <i>somatic</i>	<b>0.26 (0.13, 0.39)</b>
	STICSA S <i>cognitive</i>	0.1 (−0.03, 0.23)
Mean safety probability	STICSA S <i>somatic</i>	<b>−0.21 (−0.34, −0.08)</b>
	STICSA T <i>cognitive</i>	<b>0.18 (0.04, 0.31)</b>
	STICSA T <i>somatic</i>	<b>−0.27 (−0.4, −0.14)</b>
Mean safety probability	STICSA S <i>cognitive</i>	−0.1 (−0.23, 0.03)
	STICSA S <i>somatic</i>	<b>0.23 (0.1, 0.37)</b>
	STICSA T <i>cognitive</i>	<b>−0.19 (−0.32, −0.05)</b>
	STICSA T <i>somatic</i>	<b>0.27 (0.13, 0.39)</b>

Effects with HPDIs excluding zero are shown in bold. IUS Intolerance of Uncertainty Scale, STICSA S State–trait Inventory of Cognitive and Somatic Anxiety, State Measure, STICSA T State-trait Inventory of Cognitive and Somatic Anxiety, Trait Measure, HPDI highest posterior density estimate.



**Fig. 5 Results of trait anxiety subscale models.** Points indicate the mean of the posterior distribution for the regression coefficient parameter, while error bars represent the 95% highest posterior density interval. The  $\beta$  estimate here refers to the regression coefficient for each predictor, estimated from 400 subjects.

estimate these effects precisely, making it unlikely that they are simply a product of statistical noise. One explanation for the discrepancy between our results and those found by Aylward et al.<sup>12</sup> is that this previous study included subjects with a mix of anxiety and depressive disorders, and a negative bias in learning

may be more characteristic of depressive symptoms. This is supported by our subscale analysis where cognitive symptoms, which align more with those of depression, were associated with faster learning from danger. Our PLS analysis provides further support for this speculation as we found that symptoms of



**Table 2** Estimates from regression model predicting learning-related variables derived from our computational model from the three transdiagnostic factors identified by Gillan et al.<sup>6</sup>.

Target variable	Predictor	Estimate ( $\pm 95\%$ 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)

Effects with HPDIs excluding zero are shown in bold.  
 AD anxious-depression, CBIT compulsive behaviour and intrusive thought, SW social withdrawal, HPDI highest posterior density estimate.

depression were associated with elevated learning from threat, suggesting that such a bias in learning is associated more with depressive symptoms. It is also possible that the nature of our game-based online task engaged processes distinct from that of standard lab-based tasks. However, we believe this is unlikely since we replicate behavioural patterns shown in more traditional tasks, and also observed similar associations with anxiety in a previous lab-based study<sup>15</sup>. As such, we are confident that this is not simply due to the task used and, taken together with this prior work, our results are suggestive of a dissociation between physiological anxiety and cognitive symptoms of anxiety and depression, where a negative bias is more characteristic of cognitive symptoms.

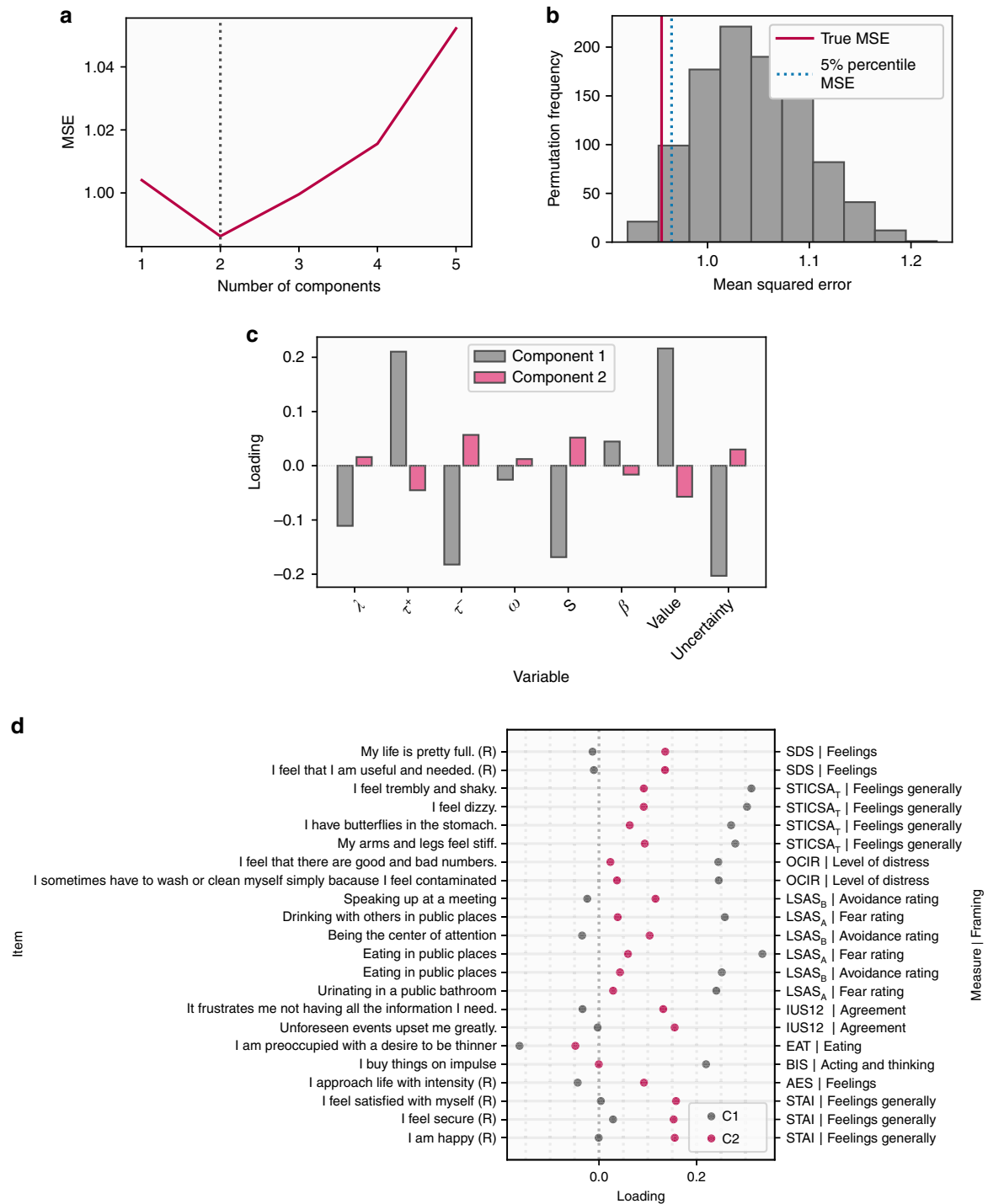
We found a similar pattern of enhanced learning from safety when examining a transdiagnostic factor representing compulsivity and intrusive thought. Although this factor has been shown to be associated with less model-based behaviour<sup>6,25</sup>, altered confidence judgements<sup>26</sup>, and action-confidence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

In addition to our a priori specified tests, we used a data-driven approach to derive components of psychopathology grounded in computational analyses of aversive learning, with the intention of providing an hypothesis-generating analysis. Using PLS regression, we identified two latent components, one broadly associating greater learning from safety with physiological symptoms of anxiety and compulsivity, while the other associated greater learning from threat with depressive symptoms and social anxiety. This method represents a conceptually similar approach to factor analytic methods used in previous large-scale online studies<sup>6</sup>, but builds upon this work by incorporating behaviour into

the process. As such, it is not surprising we identify qualitatively different components to prior work<sup>6</sup> based purely on self-report data. Notably, this data-driven analysis also revealed relationships between aversive learning and impulsive behaviour, encompassing a symptom dimension that is typically studied in the context of reward processing<sup>40</sup>. Individuals scoring higher on these symptoms exhibited higher safety learning, which may explain previously observed relationships between impulsivity and risk tolerance<sup>41</sup>. While this analysis was exploratory, we demonstrate its robustness through testing on held-out data our results are not affected by overfitting<sup>37</sup>, and these results broadly reflect those evident in our anxiety subscale analyses where physiological anxiety was associated with enhanced safety learning while cognitive anxiety was associated with enhanced threat learning. These results should be interpreted with caution due to their exploratory nature, and it will also be important to test whether this factorisation replicates across other tasks and samples.

Overall, the present results add to the growing literature showing associations between psychopathology and learning under uncertainty. Previous studies using computational approaches have largely focused on learning about rewards and losses<sup>10–12,27,42</sup>, or perceptual learning<sup>9</sup>, and those that have used more aversive paradigms (using outcomes intended to evoke subjective anxiety), such as learning to predict electric shocks, have been limited by small samples<sup>5,15,18,43</sup>. While there is a rich literature using simple fear conditioning paradigms to investigate aversive learning in individuals with anxiety disorders<sup>44,45</sup>, these tasks typically do not manipulate uncertainty, as was the intention in our task. As a result, the precise role played by aversive learning processes in psychiatric symptoms has been unclear. Our work adds to this literature by providing an account of how these processes relate to symptoms across a range of traditional diagnostic categories, showing a dissociation between cognitive and physiological symptoms of anxiety. In addition, this study adds to the growing field of computational psychiatry, supporting proposals that investigating the computational basis of learning and decision-making can provide insights into psychiatric disorders<sup>46</sup>.

A further important feature of this study is our development of an online task for measuring aversive learning. A number of studies examining other aspects of learning and decision-making in the context of psychiatric disorders have also availed of large samples recruited through online services<sup>6,8,25,26</sup>. However, it has been difficult to examine aversive learning in online environments, as aversive lab stimuli such as shock cannot be easily administered online. Only one study thus far has investigated threat-related decision-making (although not learning) online, using monetary loss as an aversive stimulus<sup>27</sup>. A game-based design allowed us to design a task that required avoidance behaviour as well as evoke feelings of anxiety, taking advantage of the well-known ability of games to produce strong emotional reactions<sup>47–51</sup>. While this is not a typical approach to designing aversive stimuli, we believe that the threat of losing points in the game provides for an aversive context. Prior work using a similar game has demonstrated that subjects show negative emotional responses to losing points<sup>52</sup>, and subjects in our task largely reported some degree of subjective anxiety and motivation to avoid the asteroids. Together, we believe this suggests that such a task can be construed as aversive, even if it does not use typical aversive stimuli such as electric shocks. In addition, although qualitatively different from standard lab-based tasks, we observed similar patterns of biased learning to that seen in lab-based work<sup>15</sup>. An added benefit of our task is that it is highly engaging, and subjects reported feeling motivated to perform well. As a result, we did not have to apply strict exclusion criteria. These features are not only important for the kind of large-scale online testing performed here. This task renders it feasible to measure



**Fig. 6 Results of PLS regression analysis.** **a** The optimal number of components was determined based on cross-validated predictive accuracy within 75% of the data used for training the model. This figure represents the mean squared error of these predictions across models with between one and five factors, showing best performance with two components. **b** Null distribution of predictive accuracy scores generated by retraining our PLS regression model on 1000 permuted datasets and testing on the held out 25% of the data set (100 subjects), with the mean squared error (MSE) achieved by the model trained on the true data shown by the red line. **c** Loadings for the two components on behavioural variables, including all parameters in the model and mean safety probably and uncertainty estimates, across all subjects. **d** Loadings on questionnaire items showing the largest dissociations in loadings between the two components, identified by taking the lowest and highest 10% of differences between loadings. Items marked (R) are reverse coded. The labels on the right represent the measure the item is taken from and an indicator of how the question is framed. SDS Zhung Self-Rating Depression Scale, STICSA State Trait Inventory of Cognitive and Somatic Anxiety, OCIR Obsessive Compulsive Inventory, LSAS Liebowitz Social Anxiety Scale (**a**, **b** represent subscales), IUS12 Intolerance of Uncertainty Scale, EAT Eating Attitudes Test, AES Apathy Evaluation Scale, STAI State Trait Anxiety Inventory.

aversive learning at regular intervals without subjects needing to physically visit the lab, a feature that could be of considerable utility in clinical trials, although such longitudinal work will first require test–retest validity of the task to be demonstrated.

With regard to computational modelling, we elected to focus our analysis on probabilistic models as these were successful in previous work, possibly reflecting their natural representation of uncertainty. However, we acknowledge that there may be alternative behavioural models that would explain the data equally well, an interpretation supported by the fact that a standard Rescorla–Wagner model performed reasonably well in our model comparison. We emphasise that our intention was not to identify the precise mechanism through which subjects are learning in this task, but instead to use a previously validated model to examine relationships with psychopathology. The computational mechanisms underlying aversive learning in humans has received a relative lack of attention compared to a more concerted focus on the reward domain.

One potential limitation of this study is a focus on a general population sample which, being recruited online, was not subject to the kind of detailed assessment possible offline. While this might limit applicability to clinical anxiety, other research indicates that findings from clinical samples replicate in samples recruited online<sup>6,8</sup>. Furthermore, it is increasingly recognised that clinical disorders lie on a continuum from health to disorder<sup>53</sup>. Although we did not deliberately set out to recruit individuals with clinically significant anxiety, 36% of our sample scored at or above a threshold designed for the detection of anxiety disorders on our measure of trait anxiety (see Supplementary Fig. 7). In light of this, and given limitations with research in clinical samples that includes medication load<sup>54</sup> and recruitment challenges<sup>55</sup>, online samples provide an effective method for studying clinically relevant phenomena. Nevertheless, replicating these findings in clinical samples is an important next step, and will be critical if we are to eventually translate results such as these to the clinic. Relatedly, we acknowledge that online samples may not necessarily be representative of the general population. On the other hand, it has been argued that online recruitment provides a more representative and inclusive sample than traditional campus-based studies<sup>56</sup>. In addition to confirming replicability across different populations, future work will need to evaluate the test–retest reliability of the task, while accounting for potential state-dependence of these measures, as suggested by our results relating to state anxiety.

In addition, it is important to note that the effects we observed were small, as in previous studies using large-scale online testing<sup>6,25</sup>. However, large samples provide accurate effect size estimates in contrast to the exaggerated effects that are common in studies using small samples<sup>57</sup>. Such small effects are unsurprising given the multifactorial nature of psychiatric disorders<sup>58</sup>. While we have shown aversive learning to be important, we acknowledge this is likely to be one of a multitude of processes involved in the development of these conditions. Finally, there was a degree of negative correlation between our update parameters for threat and safety, which may have limited our power to detect relationships between these measures and symptoms.

The results we report suggest several directions for future research. In particular, the finding that more physiologically anxious individuals tend to overestimate safety likelihood runs counter to intuition, and further work is required to understand how this may relate to symptom expression. One speculative possibility is that a persistent underestimation of threat likelihood would lead to an abundance of aversive prediction errors, causing a state of subjective physiological anxiety. An alternative explanation is the result reflects a tendency for highly anxious individuals to seek safety, and be resistant to leaving places associated

with safety<sup>59,60</sup>. However, these hypotheses await direct testing, and it will be especially important to examine them in large-scale clinical samples, taking into account a broader range of psychiatric phenotypes. Another important aspect of learning uncertainty that we did not investigate is volatility, namely the tendency of stimulus–outcome relationships to change over time. There is an evidence that individuals high in trait anxiety fail to adapt to volatility, and this deserves further study in relation to transdiagnostic psychiatric symptoms<sup>5</sup>. Furthermore, future research should investigate how aversive learning relates to psychiatric symptoms within subject. Some of our strongest results relate to state rather than trait anxiety, suggesting that aversive learning processes might fluctuate within subject depending on their affective and physiological state.

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

## Methods

**Ethics.** This research was approved by the University College London research ethics committee (reference 9929/003) and complied with all relevant ethical regulations. All participants provided informed consent by completing a form including checkboxes that were clicked to indicate agreement with various aspects of the study and overall consent to participate, and were compensated financially for their time at a rate of at least £6 per hour.

**Participants.** We recruited 400 participants through Prolific<sup>33</sup>. Subjects were selected based on being aged 18–65 and having at least a 90% approval rate across studies they had previously participated in. As described in our preregistration, we used a precision-based stopping rule to determine our sample size, stopping at the point at which either the 95% highest posterior density interval (HPDI) for all effects in our regression model reached 0.15 (checking with each 50 subjects recruited) or we had recruited 400 subjects. The precision target was not reached, and so we stopped at 400 subjects.

**Avoidance learning task.** Traditional lab-based threat learning tasks typically use aversive stimuli such as electric shocks as outcomes to be avoided. As it is not possible to use these stimuli online, we developed a game-based task in which subjects' goal was to avoid negative outcomes. While no primary aversive stimuli were used, and subjects received no actual monetary reward, there is an extensive literature showing that video games without such outcomes evoke strong positive and negative emotional experiences<sup>47–51</sup>, making this a promising method for designing an aversive learning task. In this game, participants were tasked with flying a spaceship through asteroid belts. Subjects were able to move the spaceship in the Y-axis alone, and this resulted in a one dimensional behavioural output. Crashing into asteroids diminished the spaceship's integrity by 10%. The spaceship's integrity slowly increased over the course of the task, however, if enough asteroids were hit the integrity reduced to zero and the game finished. In this eventuality subjects were able to restart and continue where they left off. The overarching goal was to maximise the number of points scored, where the latter accumulated continuously for as long as the game was ongoing, and reset if the spaceship was destroyed. Subjects were shown the current integrity of the spaceship by a bar displayed in the corner of the screen, along with by a display of their current score.

Crucially, the location of safe spaces in the asteroid belts could be learned, and learning facilitated performance as it allowed correct positioning of the spaceship prior to observing the safe location. The task was designed such that without such pre-emptive positioning it was near impossible to successfully avoid the asteroids, thus encouraging subjects to learn the safest positions. Holes in the asteroids could appear either at the top or bottom of the screen (Fig. 1a), and the probability of safety associated with either location varied independently over the course of the task. Thus, it was possible to learn the safety probability associated with each safety zone and adapt one's behaviour accordingly. The probability of each zone being safe was largely independent from the other (so that observing safety in one zone did not necessarily indicate the other was dangerous), although at least one zone was always safe on each trial. This was important, because if outcomes were entirely symmetric (i.e. safety in one zone indicated danger in the other), we would be unable to determine the extent to which value updating was driven by safety versus danger. Thus, our task aimed to largely dissociate learning from threat and safety, as outcomes are not entirely symmetric.

Trials were designed such that one option had a 90% or 10% chance of being safe for the duration of between 20 and 80 trials, subject to the condition that on a

particular trial one zone had to be safe (so that the subject had an opportunity to avoid the asteroids). This design feature meant that at any time either both, or just one, zone had a high safety likelihood. Safety probability was designed to fluctuate relatively rapidly to ensure that uncertainty fluctuated continuously over the course of the task. The probabilistic nature of the task ensured that behaviour was not straightforwardly dictated by the outcomes observed in the task. For example, encountering danger did not necessarily mean subjects should change their position on the following trial, as this outcome could be a chance event rather than signifying the chosen zone is no longer the safest.

Participants also completed a control task that required avoidance that was not dependent on learning, enabling us to control for general motor-related avoidance ability in further analyses (described in Supplementary Methods). After completing the task, subjects were asked to provide ratings indicating how anxious the task made them feel and how motivated they were to avoid the asteroids, using visual analogue scales ranging from 0 to 100.

**Inclusion and exclusion criteria.** We included subjects based on their age (18–65 years) and having a 90% prior approval rate on Prolific. We elected a priori to exclude subjects with limited response variability (indicated by a standard deviation of their positions below 0.05) so as to remove subjects who did not move the spaceship, and subjects who had missing data. However, no subject met these exclusion criteria.

**Behavioural data extraction.** For analysis, we treated each pass through an asteroid belt as a trial. Overall there were 269 trials in total. As a measure of behaviour, we extracted the mean  $Y$  position across the 1 s prior to observing the asteroid belt, representing where subjects were positioning themselves in preparation for the upcoming asteroid belt. This  $Y$  position was used for subsequent model fitting. On each trial, the outcome for each zone was regarded as danger if asteroids were observed (regardless of whether they were hit by the subject) or safety if a hole in the asteroid belt was observed.

**Computational modelling of behaviour.** Our modelling approach focused on models that allowed the quantification of subjective uncertainty. To this end, we modelled behaviour using approximate Bayesian models that assume subjects estimate safety probability using a beta distribution. This approach is naturally suited to probability estimation tasks, as the beta distribution is bounded between zero and one, and provides a measure of uncertainty through the variance of the distribution. While certain reinforcement learning formulations can achieve similar uncertainty-dependent learning and quantification of uncertainty, we chose beta models as they have an advantage of being computationally simple. Empirically, these models have been used successfully in previous studies to capture value-based learning<sup>61</sup>, where they explain behaviour in aversive learning tasks better than commonly used reinforcement learning models<sup>15,62</sup>, a pertinent characteristic in the current task.

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

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

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

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

We also observed in pilot data that subjects tended to be influenced more by outcomes occurring in the zone they had previously chosen, an effect likely due to

attention. On this basis, we incorporated a weighting parameter that allowed the outcome of the unchosen option to be downweighted by an amount shown in the above equation ( $W$ ) determined by an additional free parameter,  $\omega$ .

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

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

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

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

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

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

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

Further models elaborated on this basic premise, and full details are provided in Supplementary Methods. For completeness, we also tested two reinforcement learning models, a Rescorla–Wagner model and a variant of this model with different learning rates for better and worse than expected outcomes<sup>63</sup>, both of which are described in Supplementary material. However, we focus on the probabilistic models due to their ability to represent uncertainty naturally; our primary aim was not to differentiate between probabilistic and reinforcement learning models, but to use previously validated models to provide insights into the relationship between aversive learning, uncertainty, and psychopathology.

Models were fit with a hierarchical Bayesian approach using variational inference implemented in PyMC3, through maximising the likelihood of the data given a reparametrised beta distribution with a mean provided by the model and a single free variance parameter. Model fit was assessed using the Watanabe–Akaike Information Criterion (WAIC)<sup>64</sup>, an index of model fit designed for Bayesian models that accounts for model complexity. Parameter distributions were visualised using raincloud plots<sup>65</sup>.

**Measures of psychiatric symptoms.** Our first set of hypotheses focused on state/trait anxiety and intolerance of uncertainty. These were measured using the State Trait Inventory of Cognitive and Somatic Anxiety (STICSA)<sup>66</sup> and the Intolerance of Uncertainty Scale (IUS)<sup>67</sup> respectively. We also wished to examine how behaviour in our task related to the three transdiagnostic factors identified by Gillan et al.<sup>6</sup>, based on factor analysis of a range of psychiatric measures (Supplementary Table 1). To measure these factors more efficiently, we developed a reduced set of questions that provided an accurate approximation of the true factor scores (Supplementary Fig. 3), details of which are provided in Supplementary Methods. We also performed analyses using an approximation of clinical anxiety status, results of which are shown in Supplementary Fig. 8.

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

These models were constructed using Bambi<sup>68</sup> and fit using Markov chain Monte Carlo sampling, each with 8000 samples, 2000 of which were used for burn-in. All models included age and sex as covariates, along with performance on our control task to account for non-learning-related avoidance ability. For analyses predicting state and trait anxiety and intolerance of uncertainty, we constructed a separate model for each variable due to the high collinearity between these measures. For analyses including the three transdiagnostic factors, these were entered into a single model. When reporting regression coefficients, we report the mean of the posterior distribution along with the 95% HPDI, representing the points between which 95% of the posterior distribution's density lies. All analyses were specified in our preregistration. We did not correct for multiple comparisons in these analyses as our approach uses Bayesian parameter estimation, rather than



frequentist null hypothesis significance testing, and as such multiple comparison correction is unnecessary and incompatible with this method<sup>69</sup>.

**PLS regression.** To provide a data-driven characterisation of the relationship between task behaviour and psychiatric symptoms, and identify transdiagnostic components that are grounded in both self-report and behaviour, we used PLS regression to identify dimensions of covariance between individual questions and the measures derived from our modelling. We excluded the STICSA state subscale from this analysis, so that only trait measures were included. To ensure robustness of these results, we split our data into training and testing sets, made up of 75 and 25% of the data, respectively. To identify the appropriate number of components within the training set, we used a tenfold cross-validation procedure, fitting the model on 90% of the training data and evaluating its performance on the left-out 10%. The mean squared error of the model's predictions was then averaged across test folds to provide an index of the model's predictive accuracy with different numbers of components, using cross-validation to reduce the risk of overfitting.

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

Recent work has highlighted the risks inherent in PLS-like methods when used in high dimensional datasets<sup>37</sup>, namely that they can easily be overfit resulting in solutions that do not generalise beyond the data used to fit the model. Our approach avoids these problems by evaluating the performance on our model 25% of the data that has been held out from the model fitting stage.

**Preregistration.** The main hypotheses and methods of this study were pre-registered on the Open Science Framework [<https://osf.io/jp5qn>]. The data-driven PLS regression analysis was exploratory.

**Statistics and reproducibility.** The reported results derive from a single experiment that was not replicated.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

All raw data supporting the findings of this study are available through the Open Science Framework at <https://osf.io/b95w2/> with the DOI <https://doi.org/10.17605/OSF.IO/B95W2>. A reporting summary for this article is available as a Supplementary Information File. Source data are provided with this paper.

## Code availability

Code is available at <https://github.com/tobywise/online-aversive-learning>. The task was programmed in Javascript using Phaser 3. Analysis was performed in Python 2.7 using PyMC3 (version 3.5), Scikit-Learn (version 0.30) and Bambi (version 0.1.5). Source data are provided with this paper.

Received: 29 April 2020; Accepted: 13 July 2020;

Published online: 21 August 2020

## References

- Reininghaus, U. et al. Stress sensitivity, aberrant salience, and threat anticipation in early psychosis: an experience sampling study. *Schizophr. Bull.* **42**, 712–722 (2016).
- Harrison, A., Tchanturia, K. & Treasure, J. Attentional bias, emotion recognition, and emotion regulation in anorexia: state or trait? *Biol. Psychiatry* **68**, 755–761 (2010).
- Friston, K. J., Stephan, K. E., Montague, R. & Dolan, R. J. Computational psychiatry: the brain as a phantastic organ. *Lancet Psychiatry* **1**, 148–158 (2014).
- Montague, P. R., Dolan, R. J., Friston, K. J. & Dayan, P. Computational psychiatry. *Trends Cogn. Sci.* **16**, 72–80 (2012).
- Browning, M., Behrens, T. E., Jocham, G., O'Reilly, J. X. & Bishop, S. J. Anxious individuals have difficulty learning the causal statistics of aversive environments. *Nat. Neurosci.* **18**, 590–596 (2015).
- Gillan, C. M., Kosinski, M., Whelan, R., Phelps, E. A. & Daw, N. D. Characterizing a psychiatric symptom dimension related to deficits in goal-directed control. *eLife* **5**, e11305 (2016).
- Huang, H., Thompson, W. & Paulus, M. P. Computational dysfunctions in anxiety: failure to differentiate signal from noise. *Biol. Psychiatry* **82**, 440–446 (2017).
- Seow, T. X. F. & Gillan, C. M. Transdiagnostic phenotyping reveals a host of metacognitive deficits implicated in compulsivity. *Sci. Rep.* **10**, 2883 (2020).
- Lawson, R. P., Mathys, C. & Rees, G. Adults with autism overestimate the volatility of the sensory environment. *Nat. Neurosci.* **20**, 1293–1299 (2017).
- Mkrtchian, A., Aylward, J., Dayan, P., Roiser, J. P. & Robinson, O. J. Modeling avoidance in mood and anxiety disorders using reinforcement learning. *Biol. Psychiatry* **82**, 532–539 (2017).
- Vinckier, F. et al. Confidence and psychosis: a neuro-computational account of contingency learning disruption by NMDA blockade. *Mol. Psychiatry* **21**, 946–955 (2016).
- Aylward, J. et al. Altered learning under uncertainty in unmedicated mood and anxiety disorders. *Nat. Hum. Behav.* **1**, <https://doi.org/10.1038/s41562-019-0628-0> (2019).
- Robinson, O. J. & Chase, H. W. Learning and choice in mood disorders: searching for the computational parameters of anhedonia. *Comput. Psychiatry*, 1–26, [https://doi.org/10.1162/CPSY\\_a\\_00009](https://doi.org/10.1162/CPSY_a_00009) (2017).
- Sharp, P. B. & Eldar, E. Computational models of anxiety: nascent efforts and future directions. *Curr. Dir. Psychol. Sci.* **28**, 170–176 (2019).
- Wise, T., Michely, J., Dayan, P. & Dolan, R. J. A computational account of threat-related attentional bias. *PLoS Comput. Biol.* **15**, e1007341 (2019).
- Pulcu, E. & Browning, M. The misestimation of uncertainty in affective disorders. *Trends in Cogn. Sci.* <https://doi.org/10.1016/j.tics.2019.07.007> (2019).
- Bach, D. R. & Dolan, R. J. Knowing how much you don't know: a neural organization of uncertainty estimates. *Nat. Rev. Neurosci.* **13**, 572–586 (2012).
- de Berker, A. O. et al. Computations of uncertainty mediate acute stress responses in humans. *Nat. Commun.* **7**, 10996 (2016).
- Mathys, C., Daunizeau, J., Friston, K. J. & Stephan, K. E. A Bayesian foundation for individual learning under uncertainty. *Front. Hum. Neurosci.* **5**, 39 (2011).
- Mathys, C. et al. Uncertainty in perception and the Hierarchical Gaussian Filter. *Front. Hum. Neurosci.* **8**, 825 (2014).
- Grupe, D. W. & Nitschke, J. B. Uncertainty and anticipation in anxiety: an integrated neurobiological and psychological perspective. *Nat. Rev. Neurosci.* **14**, 488–501 (2013).
- Dugas, M. J., Gagnon, F., Ladouceur, R. & Freeston, M. H. Generalized anxiety disorder: a preliminary test of a conceptual model. *Behav. Res. Ther.* **36**, 215–226 (1998).
- Barlow, D. H. Unraveling the mysteries of anxiety and its disorders from the perspective of emotion theory. *Am. Psychol.* **55**, 1247–1263 (2000).
- Grillon, C. Associative learning deficits increase symptoms of anxiety in humans. *Biol. Psychiatry* **51**, 851–858 (2002).
- Patzelt, E. H., Kool, W., Millner, A. J. & Gershman, S. J. Incentives boost model-based control across a range of severity on several psychiatric constructs. *Biol. Psychiatry*, <https://doi.org/10.1016/j.biopsych.2018.06.018> (2018).
- Rouault, M., Seow, T., Gillan, C. M. & Fleming, S. M. Psychiatric symptom dimensions are associated with dissociable shifts in metacognition but not task performance. *Biol. Psychiatry* **84**, 443–451 (2018).
- Norbury, A., Robbins, T. W. & Seymour, B. Value generalization in human avoidance learning. *eLife* **7**, e34779 (2018).
- Michelini, G. et al. Delineating and validating higher-order dimensions of psychopathology in the Adolescent Brain Cognitive Development (ABCD) study. *Transl. Psychiatry* **9**, 1–15 (2019).
- Blanco, C. et al. Mapping common psychiatric disorders: structure and predictive validity in the National Epidemiologic Survey on Alcohol and Related Conditions. *JAMA Psychiatry* **70**, 199–207 (2013).
- Røysamb, E. et al. The joint structure of DSM-IV Axis I and Axis II disorders. *J. Abnorm. Psychol.* **120**, 198–209 (2011).
- Daniel-Watanabe, L., McLaughlin, M., Gormley, S. & Robinson, O. J. Association between a directly translated cognitive measure of negative bias and self-reported psychiatric symptoms. *Biol. Psychiatry: Cogn. Neurosci. Neuroimaging*, <https://doi.org/10.1016/j.bpsc.2020.02.010> (2020).
- Vaghi, M. M. et al. Compulsivity reveals a novel dissociation between action and confidence. *Neuron* **96**, 348–354 (2017).
- Palan, S. & Schitter, C. Prolific.ac—a subject pool for online experiments. *J. Behav. Exp. Financ.* **17**, 22–27 (2018).
- Kebets, V. et al. Somatosensory-motor dysconnectivity spans multiple transdiagnostic dimensions of psychopathology. *Biol. Psychiatry* **86**, 779–791 (2019).
- Jessen, K. et al. Patterns of cortical structures and cognition in antipsychotic-naïve patients with first-episode schizophrenia: a partial least squares correlation analysis. *Biol. Psychiatry: Cogn. Neurosci. Neuroimaging* **4**, 444–453 (2019).
- Mihalik, A. et al. Multiple hold-outs with stability: improving the generalizability of machine learning analyses of brain-behaviour relationships:

- a novel framework to link behaviour to neurobiology. *Biol. Psychiatry*, <https://doi.org/10.1016/j.biopsych.2019.12.001> (2019).
37. Dinga, R. et al. Evaluating the evidence for biotypes of depression: Methodological replication and extension of Drysdale et al. (2017). *NeuroImage Clin.* **22**, 101796 (2019).
  38. Wakeling, I. N. & Morris, J. J. A test of significance for partial least squares regression. *J. Chemom.* **7**, 291–304 (1993).
  39. Mevik, B.-H. & Cederkvist, H. R. Mean squared error of prediction (MSEP) estimates for principal component regression (PCR) and partial least squares regression (PLSR). *J. Chemom.* **18**, 422–429 (2004).
  40. Plichta, M. M. & Scheres, A. Ventral-striatal responsiveness during reward anticipation in ADHD and its relation to trait impulsivity in the healthy population: a meta-analytic review of the fMRI literature. *Neurosci. Biobehav. Rev.* **38**, 125–134 (2014).
  41. Deakin, J., Aitken, M., Robbins, T. & Sahakian, B. J. Risk taking during decision-making in normal volunteers changes with age. *J. Int. Neuropsychol. Soc.* **10**, 590–598 (2004).
  42. Pulcu, E. & Browning, M. Affective bias as a rational response to the statistics of rewards and punishments. *eLife Sci.* **6**, e27879 (2017).
  43. Michely, J., Rigoli, F., Rutledge, R. B., Hauser, T. U. & Dolan, R. J. Distinct processing of aversive experience in amygdala subregions. *Biol. Psychiatry. Cogn. Neurosci. Neuroimaging* <https://doi.org/10.1016/j.bpsc.2019.07.008> (2019).
  44. Lissek, S. et al. Classical fear conditioning in the anxiety disorders: a meta-analysis. *Behav. Res. Ther.* **43**, 1391–1424 (2005).
  45. Duits, P. et al. Updated meta-analysis of classical fear conditioning in the anxiety disorders. *Depression Anxiety* **32**, 239–253 (2015).
  46. Browning, M. et al. Realizing the clinical potential of computational psychiatry: report From the Banbury Center Meeting, February 2019. *Biol. Psychiatry*, <https://doi.org/10.1016/j.biopsych.2019.12.026> (2020).
  47. Schneider, E. F., Lang, A., Shin, M. & Bradley, S. D. Death with A Story: How Story Impacts Emotional, Motivational, and Physiological Responses to First-person Shooter Video Games. *Hum. Commun. Res.* **30**, 361–375 (2004).
  48. Ravaja, N., Saari, T., Salminen, M., Laarni, J. & Kallinen, K. Phasic emotional reactions to video game events: a psychophysiological investigation. *Media Psychol.* **8**, 343–367 (2006).
  49. Ravaja, N., Turpeinen, M., Saari, T., Puttonen, S. & Keltikangas-Järvinen, L. The psychophysiology of James Bond: phasic emotional responses to violent video game events. *Emotion* **8**, 114–120 (2008).
  50. Hutton, E. & Sundar, S. S. Can video games enhance creativity? Effects of emotion generated by dance dance revolution. *Creativity Res. J.* **22**, 294–303 (2010).
  51. Ravaja, N. et al. Spatial presence and emotions during video game playing: does it matter with whom you play? *Presence: Teleoperators Virtual Environ.* **15**, 381–392 (2006).
  52. Johnstone, T., van Reekum, C. M., Hird, K., Kirsner, K. & Scherer, K. R. Affective speech elicited with a computer game. *Emotion* **5**, 513–518 (2005).
  53. Plomin, R., Haworth, C. M. A. & Davis, O. S. P. Common disorders are quantitative traits. *Nat. Rev. Genet.* **10**, 872–878 (2009).
  54. Phillips, M. L., Travis, M. J., Fagioli, A. & Kupfer, D. J. Medication effects in neuroimaging studies of bipolar disorder. *Am. J. Psychiatry* **165**, 313–320 (2008).
  55. Wise, T. et al. Recruiting for research studies using online public advertisements: examples from research in affective disorders. *Neuropsychiatr. Dis. Treat.* **12**, 279–285 (2016).
  56. Gillan, C. M. & Daw, N. D. Taking psychiatry research online. *Neuron* **91**, 19–23 (2016).
  57. Gelman, A. & Carlin, J. Beyond power calculations: assessing type S (sign) and type M (magnitude) errors. *Perspect. Psychol. Sci.* **9**, 641–651 (2014).
  58. Kendler, K. S. From many to one to many—the search for causes of psychiatric illness. *JAMA Psychiatry* **76**, 1085–1091 (2019).
  59. Beesdo-Baum, K. et al. Avoidance, safety behavior, and reassurance seeking in generalized anxiety disorder. *Depression Anxiety* **29**, 948–957 (2012).
  60. Salkovskis, P. M. The importance of behaviour in the maintenance of anxiety and panic: a cognitive account. *Behav. Cogn. Psychother.* **19**, 6–19 (1991).
  61. de Boer, L. et al. Attenuation of dopamine-modulated prefrontal value signals underlies probabilistic reward learning deficits in old age. *eLife* **6**, e26424 (2017).
  62. Tzovara, A., Korn, C. W. & Bach, D. R. Human pavlovian fear conditioning conforms to probabilistic learning. *PLoS Comput. Biol.* **14**, e1006243 (2018).
  63. Lefebvre, G., Lebreton, M., Meyniel, F., Bourgeois-Gironde, S. & Palminteri, S. Behavioural and neural characterization of optimistic reinforcement learning. *Nat. Hum. Behav.* **1**, 0067 (2017).
  64. Watanabe, S. Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *J. Mach. Learn. Res.* **11**, 3571–3594 (2010).
  65. Allen, M., Poggiali, D., Whitaker, K., Marshall, T. R. & Kievit, R. A. Raincloud plots: a multi-platform tool for robust data visualization. *Wellcome Open Res.* **4**, 63 (2019).
  66. Grös, D. F., Antony, M. M., Simms, L. J. & McCabe, R. E. Psychometric properties of the State-Trait Inventory for Cognitive and Somatic Anxiety (STICSA): comparison to the State-Trait Anxiety Inventory (STAI). *Psychol. Assess.* **19**, 369–381 (2007).
  67. Carleton, R. N., Norton, M. A. P. J. & Asmundson, G. J. G. Fearing the unknown: a short version of the Intolerance of Uncertainty Scale. *J. Anxiety Disord.* **21**, 105–117 (2007).
  68. Yarkoni, T. & Westfall, J. Bambi: A simple interface for fitting Bayesian mixed effects models. Preprint at <https://osf.io/rv7sn> (2016).
  69. Kruschke, J. K. & Liddell, T. M. The Bayesian new statistics: hypothesis testing, estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychon. Bull. Rev.* **25**, 178–206 (2018).

### Acknowledgements

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

### Author contributions

T.W. and R.J.D. designed the study. T.W. collected and analysed the data. T.W. and R.J.D. wrote the paper.

### Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41467-020-17977-w>.

**Correspondence** and requests for materials should be addressed to T.W.

**Peer review information** *Nature Communications* thanks Martin Paulus and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020