

---

# Stochastic Approximation of Gaussian Free Energy for Risk-Sensitive Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We introduce a stochastic approximation rule for estimating the free energy from  
2 i.i.d. samples generated by a Gaussian distribution with unknown mean and vari-  
3 ance. The rule is a simple modification of the Rescorla-Wagner rule, where the  
4 (sigmoidal) stimulus is taken to be either the event of over- or underestimating a  
5 target value. Since the Gaussian free energy is known to be a certainty-equivalent  
6 sensitive to the mean and the variance, the learning rule has applications in risk-  
7 sensitive decision-making. In particular, we show how to use the rule in combina-  
8 tion with the temporal-difference error in order to obtain risk-sensitive, model-free  
9 reinforcement learning algorithms.

## 10 1 Introduction

11 **Main contribution.** Let  $N(x; \mu, \rho) = \sqrt{\frac{\rho}{2\pi}} \exp\{-\frac{\rho}{2}(x - \mu)^2\}$  be the Gaussian pdf with mean  $\mu$   
12 and precision  $\rho$ . Given a sequence  $x_1, x_2, \dots$  of i.i.d. samples drawn from  $N(x; \mu, \rho)$  with unknown  
13  $\mu$  and  $\rho$ , consider the problem of estimating the free energy  $\mathbf{F}_\beta$  for a given inverse temperature  $\beta \in \mathbb{R}$ ,  
14 that is

$$\mathbf{F}_\beta = \frac{1}{\beta} \log \int_{\mathbb{R}} N(x; \mu, \rho) \exp\{\beta x\} dx = \mu + \frac{\beta}{2\rho}. \quad (1)$$

15 This paper shows that (1) can be estimated using a surprisingly simple stochastic approximation rule.  
16 If  $v \in \mathbb{R}$  is the current estimate and a new sample  $x$  arrives, update  $v$  according to

$$v \leftarrow v + 2\alpha \cdot \sigma_\beta(x - v) \cdot (x - v), \quad (2)$$

17 where  $\alpha > 0$  is a learning rate and  $\sigma_\beta(z)$  is the scaled logistic sigmoid

$$\sigma_\beta(z) = \frac{1}{1 + \exp\{-\beta z\}}. \quad (3)$$

18 The unique and stable fixed point of the learning rule (2) is equal to the desired free energy value  
19  $v^* = \mu + \frac{\beta}{2\rho}$ .

20 **Motivation.** Risk-sensitivity, the susceptibility to the higher-order moments of the return, is neces-  
21 sary for the real-world deployment of AI agents. Wrong assumptions, lack of data, misspecification,  
22 limited computation, and adversarial attacks are just a handful of the countless sources of unforeseen  
23 perturbations that could be present at deployment time. Such perturbations can easily destabilize  
24 risk-neutral policies, because they only focus on maximizing expected return while entirely neglecting  
25 the variance. This poses serious safety concerns (Russell et al., 2015; Amodei et al., 2016; Leike  
26 et al., 2017).

27 Risk-sensitive control has a long history in control theory (Coraluppi, 1997) and is an active area  
 28 of research within reinforcement learning (RL). There are multiple different approaches to risk-  
 29 sensitivity in RL: for instance in *Minimax RL*, inspired by classical robust control theory, one derives  
 30 a conservative worst-case policy over MDP parameter intervals (Nilim and El Ghaoui, 2005; Tamar  
 31 et al., 2014); and the more recent *CVaR approach* relies on using the conditional-value-at-risk as a  
 32 robust performance measure (Galichet et al., 2013; Cassel et al., 2018). We refer the reader to García  
 33 and Fernández (2015) for a comprehensive overview. Here we focus on one of the earliest and most  
 34 popular approaches (see references), consisting of the use of exponentially-transformed values, or  
 35 equivalently, the free energy as the risk-sensitive certainty-equivalent (Bellman, 1957; Howard and  
 36 Matheson, 1972).

37 The certainty-equivalent of a stochastic value  $X \in \mathbb{R}$  is defined as the representative deterministic  
 38 value  $v \in \mathbb{R}$  that a decision-maker uses as a summary of  $X$  for valuation purposes. To illustrate,  
 39 consider a first-order Markov chain over discrete states  $\mathcal{S}$  with transition kernel  $P(s'|s)$ , state-emitted  
 40 rewards  $R(s) \in \mathbb{R}$ , and discount factor  $\gamma \in [0, 1)$ . Typically RL methods use the expectation as the  
 41 certainty-equivalent of stochastic transitions (Bertsekas and Tsitsiklis, 1995; Sutton and Barto, 2018).  
 42 Therefore they compute the value  $V(s)$  of the current state  $s \in \mathcal{S}$  by (recursively) aggregating the  
 43 future values through their expectation, e.g.

$$V(s) = \int P(s'|s)\{R(s') + \gamma V(s')\} ds'. \quad (4)$$

44 Instead, Howard and Matheson (1972) proposed using the free energy as the certainty-equivalent,  
 45 that is,

$$V(s) = F_\beta(s) = \frac{1}{\beta} \log \int P(s'|s) \exp\{\beta[R(s') + \gamma V(s')]\} ds', \quad (5)$$

46 where  $\beta \in \mathbb{R}$  is the inverse temperature parameter which determines whether the aggregation is  
 47 risk-averse ( $\beta < 0$ ), risk-seeking ( $\beta > 0$ ), or even risk-neutral as a special case ( $\beta = 0$ ). Indeed, if  
 48 the future values are bounded, then  $F_\beta(s)$  is sigmoidal in shape as a function of  $\beta$ , with three special  
 49 values given by

$$\lim_{\beta} F_\beta(s) = \begin{cases} \min_{s'}\{R(s') + \gamma V(s')\} & \text{if } \beta \rightarrow -\infty; \\ \mathbf{E}[R(S') + \gamma V(S')|S = s] & \text{if } \beta \rightarrow 0; \\ \max_{s'}\{R(s') + \gamma V(s')\} & \text{if } \beta \rightarrow +\infty. \end{cases} \quad (6)$$

50 These limit values highlight the sensitivity to the higher-order moments of the return. Because of this  
 51 property, the free energy has been used as the certainty-equivalent for assessing the value of both  
 52 actions and observations under limited control and model uncertainty respectively, each effect having  
 53 their own inverse temperature. The work by Grau-Moya et al. (2016) is a demonstration of how to  
 54 incorporate multiple types of effects in MDPs.

55 The present work addresses a longstanding problem pointed out by Mihatsch and Neuneier (2002).  
 56 An advantage of using expectations is that certainty-equivalents such as (4) are easily estimated  
 57 using stochastic approximation schemes. For instance, consider the classical Robbins-Monro update  
 58 (Robbins and Monro, 1951)

$$v \leftarrow v + \alpha \cdot (x - v) \quad (7)$$

59 where  $x \sim P(x)$  is a stochastic target value,  $\alpha$  is a learning rate, and  $v$  is the estimate of  $\mathbf{E}[X]$ .  
 60 Substituting  $x = R(s') + \gamma V(s')$  and  $v = V(s)$  leads to the popular TD(0) update (Sutton and Barto,  
 61 1990):

$$V(s) \leftarrow V(s) + \alpha(R(s') + \gamma V(s') - V(s)). \quad (8)$$

62 However, there is no model-free counterpart for estimating free energies (5) under general unknown  
 63 distributions. The difficulty lies in that model-free updates rely on single (Monte-Carlo) unbiased  
 64 samples, but these are not available in the case of the free energy due to the log-term on the r.h.s.  
 65 of (5). This shortcoming led Mihatsch and Neuneier (2002) to propose the alternative risk-sensitive  
 66 learning rule

$$v \leftarrow v + \alpha \cdot u \cdot (x - v), \quad \text{where } u = \begin{cases} (1 - \kappa) & \text{if } (x - v) \geq 0 \\ (1 + \kappa) & \text{if } (x - v) < 0 \end{cases} \quad (9)$$

67 and where  $\kappa \in [-1; 1]$  is a risk-sensitivity parameter. While the heuristic (9) does produce risk-  
 68 sensitive policies, these have no formal correspondence to free energies.

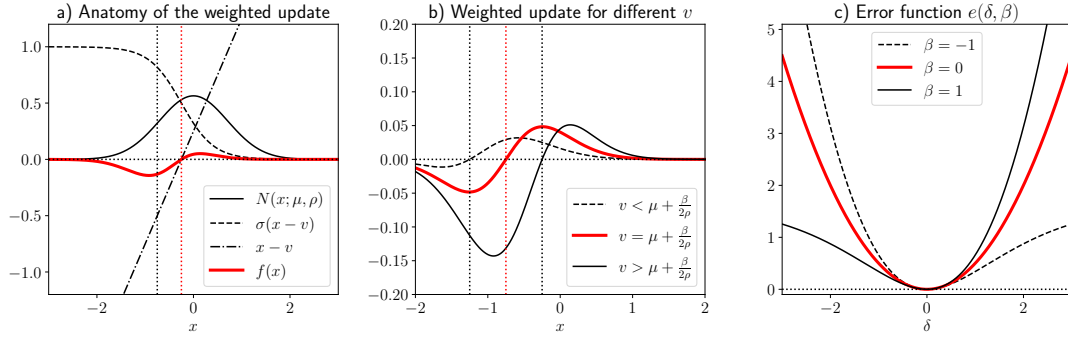


Figure 1: Update rule and its error function. **a)** shows the update  $f(x)$  to the estimate  $v$  caused by the arrival of a sample  $x$ , weighted by its probability density. The expected update is determined by comparing the integrals of the positive and negative lobes. **b)** Illustration of weighted update functions  $f(x)$  for different values of the current estimate  $v$ . The positive lobes are either larger, equal, or smaller than the negative lobes for a  $v$  that is either smaller, equal, or larger than the free energy respectively. **c)** Error function implied by the update rule. For a risk-neutral ( $\beta = 0$ ) estimator the error function is equal to the quadratic error  $e(\delta, 0) = \frac{1}{2}\delta^2$ . For a risk-averse estimator ( $\beta < 0$ ), the error function is lopsided, penalizing under-estimates stronger than over-estimates. Furthermore,  $e(\delta, \beta)$  is an even function in  $\beta$ .

69 As anticipated, our work contributes a simple model-free rule for estimating the free energy in the  
70 special case of Gaussian distributions. Starting from the Rescorla-Wagner rule

$$v \leftarrow v + \alpha \cdot u \cdot (x - v), \quad (10)$$

71 where  $u \in \{0, 1\}$  is an indicator function marking the presence of a stimulus (Rescorla, 1972), we  
72 substitute  $u$  by twice the soft-indicator function  $\sigma_\beta(x - v)$  of (3), which activates whenever  $v$  either  
73 over- or underestimates the target value  $x$ , depending on the sign of the risk-sensitivity parameter  $\beta$ .  
74 Using the substitutions appropriate for RL, we obtain the risk-sensitive TD(0)-rule

$$V(s) \leftarrow V(s) + 2\alpha \cdot \sigma_\beta(\delta) \cdot \delta, \quad (11)$$

75 where  $\delta = R(s') + \gamma V(s') - V(s)$  is the standard temporal-difference error. The learning rule is  
76 trivial to implement, works as stated for tabular RL, and is easily adapted to the objective functions  
77 of deep RL methods (Mnih et al., 2015). Finally, the learning rule is also consistent with findings in  
78 computational neuroscience (Niv et al., 2012), e.g. predicting asymmetric updates that are stronger  
79 for negative prediction errors in the risk-averse case (Gershman, 2015).

## 80 2 Analysis of the Learning Rule

81 Our central result is the following lemma, which implies that the unique and stable fixed point of the  
82 expected learning dynamics of (2) is given by the desired free energy.

83 **Lemma 1.** *If  $x_1, x_2, \dots$  are i.i.d. samples from  $P(X) = N(x; \mu, \rho)$ , then the expected update  $J(v)$*   
84 *of the learning rule (2) is twice differentiable and such that*

$$J(v) = 2\mathbf{E}[\sigma_\beta(X - v) \cdot (X - v)] \begin{cases} < 0, & \text{if } v > \mathbf{F}_\beta; \\ = 0, & \text{if } v = \mathbf{F}_\beta; \\ > 0, & \text{if } v < \mathbf{F}_\beta. \end{cases}$$

85 *Proof.* The expected update of  $v$  is

$$J(v) := 2 \int N(x; \mu, \rho) \sigma(x - v) (x - v) dx, \quad (12)$$

86 where we have dropped the subscript  $\beta$  from  $\sigma_\beta$  for simplicity. Using the Leibnitz integral rule it  
87 is easily seen that this function is twice differentiable w.r.t.  $v$ , because the integrand is a product of  
88 twice differentiable functions.

89 The resulting update direction will be positive if the integral over the positive contributions outweighs  
 90 the negative contributions and vice versa. The integrand of (12) has a symmetry property: splitting  
 91 the domain of integration  $\mathbb{R}$  into  $(-\infty; v]$  and  $(v; +\infty)$ , using the change of variable  $\delta = x - v$ , and  
 92 recombining the two integrals into one gives

$$J(v) := 2 \int_0^\infty \left\{ N(v + \delta; \mu, \rho) \sigma(\delta) - N(v - \delta; \mu, \rho) \sigma(-\delta) \right\} \delta d\delta. \quad (13)$$

93 We will show that the integrand of (13) is either negative, zero, or positive, depending on the value  
 94 of  $v$ . Define the weighted update  $f(x)$  as

$$f(x) = f(v + \delta) := N(v + \delta; \mu, \rho) \sigma(\delta) \delta.$$

95 This function is illustrated in Figure 1a. We are interested in the ratio

$$\frac{f(v + \delta)}{f(v - \delta)} = \frac{N(v + \delta; \mu, \rho) \sigma(\delta)}{N(v - \delta; \mu, \rho) \sigma(-\delta)}, \quad (14)$$

96 which compares the positive against the negative contributions to the integrand in (13). The first  
 97 fraction of the r.h.s. of (14) is equal to

$$\frac{N(v + \delta; \mu, \rho)}{N(v - \delta; \mu, \rho)} = \exp\left\{-\frac{\rho}{2}(v + \delta - \mu)^2 + \frac{\rho}{2}(v - \delta - \mu)^2\right\} = \exp\{-2\rho\delta(v - \mu)\}.$$

98 Using the symmetry property  $\sigma(\delta) = 1 - \sigma(-\delta)$  of the logistic sigmoid function, the second fraction  
 99 can be shown to be equal to

$$\frac{\sigma(\delta)}{\sigma(-\delta)} = \frac{\sigma(\delta)}{1 - \sigma(\delta)} = \exp\{\beta\delta\}.$$

100 Substituting the above back into (14) results in

$$\frac{f(v + \delta)}{f(v - \delta)} = \exp\{-2\rho\delta(v - \mu) + \beta\delta\} \begin{cases} > 1 & \text{for } v < \mu + \frac{\beta}{2\rho}, \\ = 1 & \text{for } v = \mu + \frac{\beta}{2\rho}, \\ < 1 & \text{for } v > \mu + \frac{\beta}{2\rho}, \end{cases}$$

101 also illustrated in Figure 1b. Therefore, the integrand in (13) is either positive ( $v < \mu + \frac{\beta}{2\rho}$ ), zero  
 102 ( $v = \mu + \frac{\beta}{2\rho}$ ), or negative ( $v > \mu + \frac{\beta}{2\rho}$ ), allowing to conclude the claim of the lemma.  $\square$

### 103 3 Additional Properties

104 We discuss additional properties in order to strengthen the intuition and to clarify the significance of  
 105 the learning rule; some practical implementation advice is given at the end.

106 **Associated free energy functional.** The Gaussian free energy  $F_\beta$  in (1) is formally related to the  
 107 valuation of risk-sensitive portfolios used in finance (Markowitz, 1952). It is well-known that the free  
 108 energy is the extremum of the *free energy functional*, defined as the Kullback-Leibler-regularized  
 109 expectation of  $X$ :

$$F_\beta[p(x)] := \mathbf{E}_p[X] - \frac{1}{\beta} \mathbf{KL}(p(x) \| N(x; \mu, \rho)). \quad (15)$$

110 This functional is convex in  $p$  for  $\beta < 0$  and concave for  $\beta > 0$ . Taking either the minimum (for  
 111  $\beta < 0$ ) or maximum (for  $\beta > 0$ ) w.r.t.  $p(x)$  yields

$$\mathbf{F}_\beta = \text{extr}_{p(x)} F_\beta[p(x)] = \left[ \mu + \frac{\beta}{\rho} \right] - \frac{1}{\beta} \left[ \frac{\beta^2}{2\rho} \right] = \mu + \frac{\beta}{2\rho} = \mathbf{E}[X] + \frac{\beta}{2} \mathbf{Var}[X], \quad (16)$$

112 that is, the Gaussian free energy is a linear function of  $\beta$ , where the intercept and the slope are equal  
 113 to the expectation and half of the variance of  $X$  respectively. The extremizer  $p^*(x)$  is the Gaussian

$$p^*(x) = \arg \text{extr}_{p(x)} F_\beta[p(x)] = N(x; \mu + \frac{\beta}{\rho}, \rho). \quad (17)$$

114 The above gives a precise meaning to the free energy as a certainty-equivalent. The choice of a  
 115 non-zero inverse temperature  $\beta$  reflects a distrust in the reference probability density  $N(x; \mu, \rho)$  as a  
 116 reliable model for  $X$ . Specifically, the magnitude of  $\beta$  quantifies the degree of distrust and the sign of  
 117  $\beta$  indicates whether it is an under- or overestimation. This distrust results in using the extremizer (17)  
 118 as a robust substitute for the original reference model for  $X$ .

119 **Game-theoretic interpretation.** In addition to the above, previous work (Ortega and Lee, 2014;  
 120 Eysenbach and Levine, 2019; Husain et al., 2021) has shown that the free energy functional has an  
 121 interpretation as a two-player game which characterizes its robustness properties. Following Ortega  
 122 and Lee (2014), computing the Legendre-Fenchel dual of the KL regularizer yields an equivalent  
 123 adversarial re-statement of the free energy functional (15), which for  $\beta > 0$  is given by

$$\max_{p(x)} \min_{c(x)} \left\{ \int p(x)[x - c(x)] dx + \int N(x; \mu, \rho) \exp\{\beta c(x)\} dx, \right\}, \quad (18)$$

124 where the perturbations  $c(x) \in \mathbb{R}$  are chosen by an adversary (Note: for the case  $\beta < 0$  one obtains a  
 125 Minimax problem over  $p(x)$  and  $c(x)$  rather than a Maximin). From this dual interpretation, one sees  
 126 that the distribution  $p(x)$  is chosen as if it were maximizing the expected value of  $x' = x - c(x)$ ,  
 127 the adversarially perturbed version of  $x$ . In turn, the adversary attempts to minimize  $x'$ , but at the  
 128 cost of an exponential penalty for  $c(x)$ . More precisely, given the distribution  $p(x)$ , the adversarial  
 129 best-response (ignoring constants) is

$$c^*(x) \stackrel{(a)}{=} \frac{1}{\beta} \log \frac{p(x)}{N(x; \mu, \rho)} \stackrel{(b)}{=} \frac{1}{2\beta} \left\{ \rho(x - \mu)^2 - \bar{\rho}(x - \bar{\mu})^2 + \log \frac{\bar{\rho}}{\rho} \right\} \stackrel{(c)}{=} x - \mathbf{F}_\beta, \quad (19)$$

130 where the equality (a) is true for any choice of  $p(x)$ ; (b) holds if  $p(x) = N(x; \bar{\mu}, \bar{\rho})$  for some mean  $\bar{\mu}$   
 131 and precision  $\bar{\rho}$ ; and where (c) holds if  $p(x)$  is the extremizer (17). Here we see that the adversarial  
 132 perturbations can be arbitrarily bad if  $p(x)$  is not chosen cautiously: for instance, for the (Gaussian)  
 133 Dirac delta

$$p(x) = N(x; \mu, \bar{\rho}) \xrightarrow{\bar{\rho} \rightarrow \infty} \delta(x - \mu) \quad \text{we get} \quad c^*(x) = \mathcal{O}\left(\log \frac{\bar{\rho}}{\rho}\right) \xrightarrow{\bar{\rho} \rightarrow \infty} +\infty. \quad (20)$$

134 **Error function.** Let  $\delta = x - v$  be the instantaneous difference between the sample and the estimate.  
 135 If the update rule (2) corresponds to a stochastic gradient descent step, then what is the error function?  
 136 That is, if

$$v \leftarrow v - \alpha \cdot \nabla_\delta e(\delta, \beta) = v + 2\alpha \cdot \sigma_\beta(\delta) \cdot \delta,$$

137 then what is  $e(\delta, \beta)$ ? Integrating the gradient  $\nabla_\delta e(\delta, \beta)$  with respect to  $\delta$  gives

$$e(\delta, \beta) = 2 \int \sigma(\delta) \delta d\delta = \frac{2\delta}{\beta} \log(1 + \exp\{\beta\delta\}) + \frac{2}{\beta^2} \text{li}_2(-\exp\{\beta\delta\}) + \frac{\pi^2}{6\beta^2}, \quad (21)$$

138 where  $\log(1 + \exp(z))$  is the *softplus function* (Dugas et al., 2001) and  $\text{li}_2(z)$  is *Spence's function*  
 139 (or dilogarithm) defined as

$$\text{li}_2(z) = - \int_0^z \frac{\log(1 - z)}{z} dz,$$

140 and where the constant of integration  $\frac{\pi^2}{6\beta^2}$  was chosen so that  $\lim_{\delta \rightarrow 0} e(\delta, \beta) = 0$  for all  $\beta \in \mathbb{R}$ . This  
 141 error function is illustrated in Figure 1c for a handful of values of  $\beta$ . In the limit  $\beta \rightarrow 0$ , the error  
 142 function becomes:

$$\lim_{\beta \rightarrow 0} e(\delta, \beta) = \frac{1}{2} \delta^2,$$

143 thus establishing a connection between the quadratic error and the proposed learning rule.

144 **Practical considerations.** The free energy learning rule (2) can be implemented as stated, for  
 145 instance either using constant learning rate  $\alpha > 0$  or using an adaptive learning rate  $\alpha_t > 0$  fulfilling  
 146 the Robbins-Monro conditions  $\sum_t \alpha_t > 0$  and  $\sum_t \alpha_t^2 < \infty$ .

147 A problem arises when most of the data falls within the near-zero saturated region of the sigmoid,  
 148 which can occur due to an unfortunate initialization of the estimate  $v$ . Since then  $\sigma_\beta(x - v) \approx 0$  for  
 149 most  $x$ , learning can be very slow. This problem can be mitigated using an affine transformation of the  
 150 sigmoid that guarantees a minimal rate  $\eta > 0$ , such as

$$\tilde{\sigma}_\beta(z) = \eta + (1 - 2\eta)\sigma_\beta(z), \quad (22)$$

151 which re-scales the sigmoid within the interval  $[\eta, 1 - \eta]$ . We have found this adjustment to work  
 152 well for  $|\beta| \approx 0$ , especially when it is only used during the first few iterations.

153 If one wishes to use the learning rule in combination with gradient-based optimization (as is typical  
 154 in a deep learning architecture), we do not recommend using the error function (21) directly. Rather,

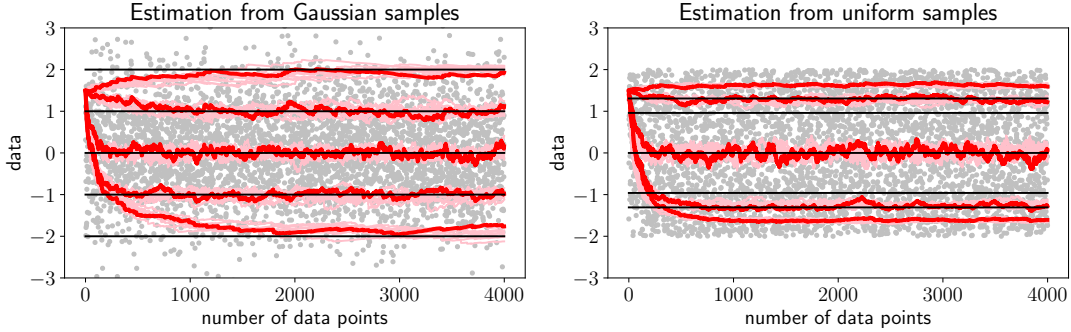


Figure 2: Estimation of the free energy from Gaussian (left panel) and uniform samples (right panel). Each plot shows 10 estimation processes (9 in pink, 1 in red) per choice of the inverse temperature, where  $\beta \in \{-4, -2, 0, 2, 4\}$ . The true free energies are shown in black. The estimation of the free energy is accurate for Gaussian data but biased for uniform data.

155 we suggest absorbing the factor  $2\tilde{\sigma}_\beta(\delta)$  directly into the learning rate (where as before,  $\delta = x - v$ ).  
 156 A simple way to achieve this consists in scaling the estimation error  $E(\delta)$  by said factor using a  
 157 stop-gradient, that is,

$$\tilde{E}(\delta) := \text{StopGrad}(2\tilde{\sigma}_\beta(\delta)) \cdot E(\delta), \quad (23)$$

158 since then the error gradient with respect to the model parameters  $\theta$  will be

$$\nabla_\theta \tilde{E}(\delta) = -2\tilde{\sigma}_\beta(\delta) \cdot \frac{\partial E}{\partial \delta} \frac{\partial v}{\partial \theta}. \quad (24)$$

159 Finally, a large  $|\beta|$  chooses a target free energy within a tail of the distribution, leading to slower  
 160 convergence. If one wishes to approximate a free energy that sits at  $n$  standard deviations from the  
 161 mean, then  $\beta$  should be chosen as

$$\beta(n) = 2n\sqrt{\rho}. \quad (25)$$

162 However, since  $\beta(n)$  is not scale invariant and the scale  $\rho$  is unknown, a good choice of  $\beta$  must be  
 163 determined empirically.

## 164 4 Experiments

165 **Estimation.** Our first experiment is a simple sanity check. We estimated the free energy in an  
 166 online manner using the learning rule (2) from data generated by two i.i.d. sources: a standard  
 167 Gaussian, and uniform distribution over the interval  $[-2, 2]$ . Five different inverse temperatures  
 168 were used ( $\beta \in \{-4, -2, 0, 2, 4\}$ ). For each condition, we ran ten estimation processes from 4000  
 169 random samples using the same starting point ( $v = 1.5$ ). The learning rate was constant and equal to  
 170  $\alpha = 0.02$ .

171 The results are shown in figure 2. In the Gaussian case, the estimation processes successfully stabilize  
 172 around the true free energies, with processes having larger  $|\beta|$  converging slower, but fluctuating  
 173 less. In the uniform case, the estimation processes do not settle around the correct free energy values  
 174 for  $\beta \neq 0$ ; however, the found solutions increase monotonically with  $\beta$ . These results validate the  
 175 estimation method only for Gaussian data, as expected.

176 **Reinforcement learning.** Next we applied the risk-sensitive learning rule to RL in a simple grid-  
 177 world. The goal was to qualitatively investigate the types of policies that result from different  
 178 risk-sensitivities. Shown in Figure 3a, the objective of the agent is to navigate to a terminal state  
 179 containing a reward pill within no more than 25 time steps while avoiding the water. The reward pill  
 180 delivers one reward point upon collection, whereas standing in the water penalizes the agent with  
 181 minus one reward point per time step. In addition, there is a very strong wind: with 50% chance in  
 182 each step, the wind pushes the agent one block in a randomly chosen cardinal direction.

183 We trained R2D2 (Kapturowski et al., 2018) agents with the risk-sensitive cost function (23) using  
 184 five uniformly spaced inverse temperatures  $\beta$  ranging from  $-0.8$  to  $0.8$ . The architecture of our

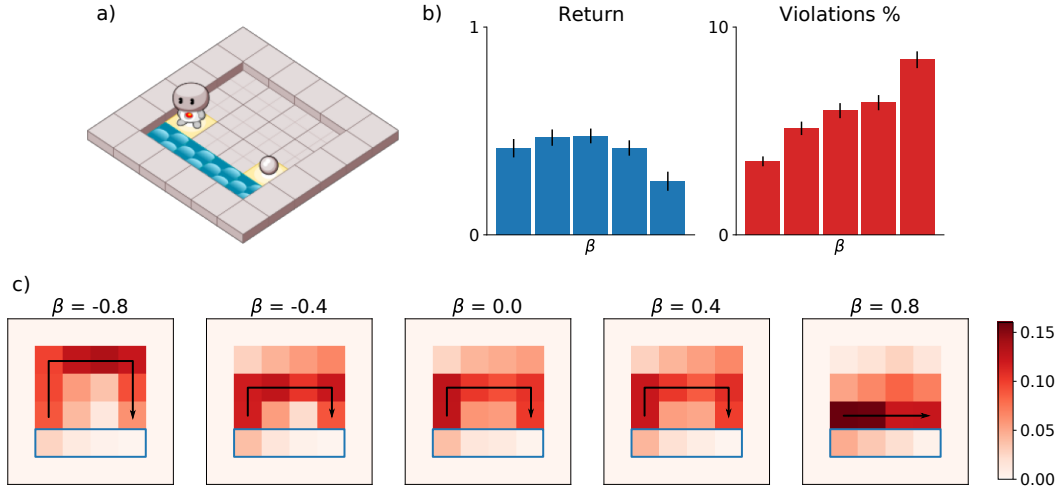


Figure 3: Comparison of risk-sensitive RL agents. **a)** The task consists in picking up a reward located at the terminal state while avoiding stepping into water. A strong wind pushes the agent into a random direction 50% of the time. **b)** Bar plots showing the average return (blue) and the percentage of violations (red) for each policy, ordered from lowest to highest  $\beta$ . **c)** State visitation frequencies for each policy, plus the optimal (deterministic) policy when there is no wind (black paths).

185 agents consisted of a first convolutional layer with 3-by-3-kernels and 128 channels, a dense layer  
 186 with 128 units, and a logit layer for the four possible actions (i.e. walking directions). The discount  
 187 factor was set to  $\gamma = 0.95$ . Each agent was trained for 500K iterations with a batch size of 64, using  
 188 the Adam optimizer with learning rate  $10^{-4}$  (Kingma and Ba, 2014). The target network was updated  
 189 every 400 steps. The inputs to the network were observation tensors of binary features representing  
 190 the 2D board. Note these agents didn’t use any recurrent cells and therefore no backpropagation  
 191 through time was used. To train all the agents in this experiment we used 154 CPU core hours at 2.4  
 192 GHz and 22.5 GPU hours.

193 To analyze the resulting policies, we computed the episodic returns and the percentage of time the  
 194 agents spent in the water (i.e. the “violations”) from 1000 roll-outs. The results, shown in Figure 3b,  
 195 reveal that the risk-neutral policy ( $\beta = 0$ ) has the highest average return. However, the percentage of  
 196 violations increases monotonically with  $\beta$ . Figure 3c shows the state-visitation probabilities estimated  
 197 from the same roll-outs. There are essentially three types of policies: risk-averse, taking the longest  
 198 path away from the water; risk-neutral, taking middle path; and risk-seeking, taking the shortest  
 199 route right next to the water. These are even more crisply revealed when the wind is de-activated.  
 200 Interestingly, the risk-averse policy ( $\beta = -0.8$ ) does not always reach the goal, which explains why  
 201 its return is slightly lower in spite of committing fewer violations.

202 **Bandits.** In the last experiment we wanted to observe the premiums that risk-sensitive agents are  
 203 willing to pay when confronted with a choice between a certain and a risky option. To do so, we  
 204 used a two-arm bandit setup, where one arm (“certain”) delivered a fixed reward and the other arm  
 205 (“risky”) a stochastic one—more precisely, drawn from a Gaussian distribution with mean  $\mu$  and  
 206 precision  $\rho = 2$ . Both the fixed payoff and the mean  $\mu$  of the risky arm were drawn from a standard  
 207 Gaussian distribution at the beginning of an episode, which lasted twenty rounds. To build agents  
 208 that can trade off exploration versus exploitation, we used memory-based meta-learning (Wang et al.,  
 209 2016; Santoro et al., 2016), which is known to produce near-optimal bandit players (Ortega et al.,  
 210 2019; Mikulik et al., 2020).

211 We meta-trained five R2D2 agents using risk-sensitives  $\beta \in \{-1.0, -0.5, 0, 0.5, 1.0\}$  on the two-  
 212 armed bandit task distribution (also randomizing the certain/risky arm positions) with discount factor  
 213  $\gamma = 0.95$ . The network architecture and training parameters were as in the previous RL experiment,  
 214 with the difference that the initial convolutional layer was replaced with a dense layer and an LSTM  
 215 layer having 128 memory cells (Hochreiter and Schmidhuber, 1997). We used backpropagation  
 216 through time for computing the episode gradients. The input to the network consisted of the action

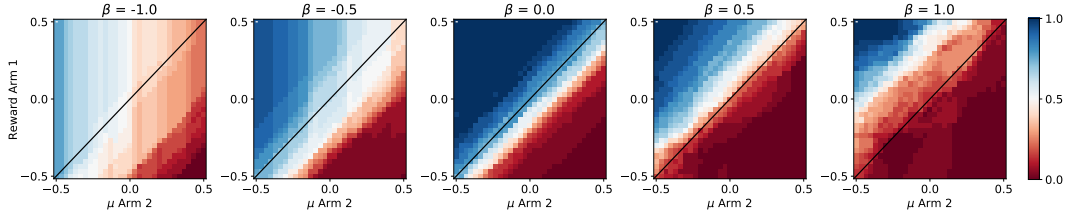


Figure 4: Two-armed bandit policy profiles with different risk-sensitivities  $\beta$ . The certain arm 1 pays a deterministic reward, while the risky arm 2 pays a stochastic reward drawn from  $N(r; \mu, \rho)$  with precision  $\rho = 2$ . The agents were meta-trained on bandits where the payoffs (i.e. arm 1’s payoff and arm 2’s mean) were drawn from a standard Gaussian distribution. The plots show the marginal probability of choosing the certain arm (blue) over the risky arm (red) after twenty interactions for every payoff combination. Each point in the uniform grid was estimated from 30 seeds. Note the deviations from the true risk-neutral indifference curve (black diagonal).

217 taken and reward obtained in the previous step. This setup allows agents to adapt their choices to past  
 218 interactions throughout an episode. To train all the agents in this experiment we used 88 CPU core  
 219 hours at 2.4 GHz and 10 GPU hours.

220 Figure 4 shows the agents’ choice profile in the last (20<sup>th</sup>) time step. A true risk-neutral agent does  
 221 not distinguish between a certain and risky option that have the same expected payoff (black diagonal).  
 222 The main finding is that the indifference region (i.e. close to a 50% choice in white color) evolves  
 223 significantly with increasing  $\beta$ , implying that the agents with different risk attitudes are indeed willing  
 224 to pay different risk premia (measured as the vertical distance of the indifference region from the  
 225 diagonal). We observe two effects. The most salient effect is that the indifference region mostly  
 226 moves from being beneath (risk-averse) to above (risk-seeking) the true risk-neutral indifference  
 227 curve as  $\beta$  increases. The second effect is that risk-averse policies ( $\beta = -1$  and  $-0.5$ ) contain a  
 228 large region of a stochastic choice profile that appears to depend only on the risky arm’s parameter.  
 229 We do not have a clear explanation for this effect. Our hypothesis is that risk-averse policies assume  
 230 adversarial environments, which require playing mixed strategies with precise probabilities. Finally,  
 231 the risk-neutral agent ( $\beta = 0$ ) appears to be slightly risk-averse. We believe that this effect arises due  
 232 to the noisy exploration policy employed during training.

## 233 5 Discussion

234 **Summary of contributions.** In this work we have introduced a learning rule for the online estima-  
 235 tion of the Gaussian free energy with unknown mean and precision/variance. The learning rule (2) is  
 236 obtained by reinterpreting the stimulus-presence indicator component of the Rescorla-Wagner rule  
 237 (Rescorla, 1972) as a (soft) indicator function for the event of either over- or underestimating the  
 238 target value. In Lemma 1 we have shown that the free energy is the unique and stable fixed point of  
 239 the expected learning dynamics. This is the main contribution.

240 Furthermore, we have shown how to use the learning rule for risk-sensitive RL. Since the free  
 241 energy implements certainty-equivalents that range from risk-averse to risk-seeking, we were able  
 242 to formulate a risk-sensitive, model-free update in the spirit of TD(0) (Sutton and Barto, 1990),  
 243 thereby addressing a longstanding problem (Mihatsch and Neuneier, 2002) for the special case of  
 244 the Gaussian distribution. Due to its simplicity, the rule is easy to incorporate into existing deep RL  
 245 algorithms, for instance by modifying the error using a stop-gradient as shown in (23). In Section 3  
 246 we also elaborated on the role of the free energy within decision-making, pointing out its robustness  
 247 properties and adversarial interpretation.

248 We also demonstrated the learning rule in experiments. Firstly, we empirically confirmed that  
 249 the online estimates stabilize around the correct Gaussian free energies (Section 4–Estimation).  
 250 Secondly, we showed how incorporating risk-attitudes into deep RL can lead to agents implementing  
 251 qualitatively different policies which intuitively make sense (Section 4–RL). Lastly, we inspected the  
 252 premia risk-sensitive agents are willing to pay for choosing a risky over a certain option, finding that  
 253 agents have choice patterns that are more complex than we had anticipated (Section 4–Bandits).



254 **Limitations.** As shown empirically in Section 4–Estimation, an important limitation of the learning  
 255 rule is that its fixed point is only equal to the free energy when the samples are Gaussian (or  
 256 approximately Gaussian, as justified by the CLT). Nevertheless, agents using the risk-sensitive TD(0)  
 257 update (11) still display risk attitudes monotonic in  $\beta$ , with  $\beta = 0$  reducing to the familiar risk-neutral  
 258 case.

259 While Lemma 1 establishes the stable equilibrium of the expected update, it only guarantees conver-  
 260 gence in continuous-time updates. To show convergence using discrete-time point samples, a stronger  
 261 result is required. In particular, we conjecture that

$$|J(v)| = 2 \left| \int N(x; \mu, \rho) \sigma_\beta(x - v)(x - v) dx \right| \leq 2|\mathbf{F}_\beta - v| \quad (26)$$

262 If (26) is true, meaning that  $J(v)$  is 2-Lipschitz, then this could be combined with a result in stochastic  
 263 approximation theory akin to Theorem 1 in Jaakkola et al. (1994) to prove convergence.

264 A shortcoming of our experiments using R2D2 agents is that they deterministically pick actions that  
 265 maximize the Q-value. However, risk-averse agents see their environments as being adversarial, and  
 266 these in turn require stochastic policies in order to achieve optimal performance.

267 **Conclusions.** Because it is impossible to anticipate the many ways in which a dynamically-changing  
 268 environment will violate prior assumptions, requiring the robustness of ML algorithms is of vital  
 269 importance for their deployment in real-world applications. Unforeseen events can render their  
 270 decisions unreliable—and in some cases even unsafe.

271 Our work makes a small but nonetheless significant contribution to risk-sensitivity in ML. In essence,  
 272 it suggests a minor modification to existing algorithms, biasing valuation estimates in a risk-sensitive  
 273 manner. In particular, we expect the risk-sensitive TD(0)-learning rule to become an integral part of  
 274 future deep RL algorithms.

## 275 References

- 276 Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. (2016). Concrete  
 277 problems in AI safety. *arXiv preprint arXiv:1606.06565*.
- 278 Bellman, R. (1957). *Dynamic Programming*. Princeton University Press.
- 279 Bertsekas, D. P. and Tsitsiklis, J. N. (1995). Neuro-dynamic programming: an overview. In  
 280 *Proceedings of 1995 34th IEEE conference on decision and control*, volume 1, pages 560–564.  
 281 IEEE.
- 282 Cassel, A., Mannor, S., and Zeevi, A. (2018). A general approach to multi-armed bandits under risk  
 283 criteria. In *Conference On Learning Theory*, pages 1295–1306. PMLR.
- 284 Coraluppi, S. P. (1997). *Optimal control of Markov decision processes for performance and robustness*.  
 285 University of Maryland, College Park.
- 286 Dugas, C., Bengio, Y., Bélisle, F., Nadeau, C., and Garcia, R. (2001). Incorporating second-order  
 287 functional knowledge for better option pricing. In *Advances in Neural Information Processing*  
 288 *Systems*, volume 13. MIT Press.
- 289 Eysenbach, B. and Levine, S. (2019). If MaxEnt RL is the answer, what is the question?
- 290 Galichet, N., Sebag, M., and Teytaud, O. (2013). Exploration vs exploitation vs safety: Risk-aware  
 291 multi-armed bandits. In *Asian Conference on Machine Learning*, pages 245–260. PMLR.
- 292 García, J. and Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. *Journal*  
 293 *of Machine Learning Research*, 16(1):1437–1480.
- 294 Gershman, S. J. (2015). Do learning rates adapt to the distribution of rewards? *Psychonomic Bulletin*  
 295 *& Review*, 22(5):1320–1327.
- 296 Grau-Moya, J., Leibfried, F., Genewein, T., and Braun, D. A. (2016). Planning with information-  
 297 processing constraints and model uncertainty in Markov decision processes. In *Joint European Con-*  
 298 *ference on Machine Learning and Knowledge Discovery in Databases*, pages 475–491. Springer.

- 299 Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–  
300 1780.
- 301 Howard, R. A. and Matheson, J. E. (1972). Risk-sensitive Markov decision processes. *Management*  
302 *science*, 18(7):356–369.
- 303 Husain, H., Ciosek, K., and Tomioka, R. (2021). Regularized policies are reward robust. In  
304 *International Conference on Artificial Intelligence and Statistics*, pages 64–72. PMLR.
- 305 Jaakkola, T., Jordan, M. I., and Singh, S. P. (1994). On the convergence of stochastic iterative  
306 dynamic programming algorithms. *Neural computation*, 6(6):1185–1201.
- 307 Kappen, H. J. (2005). Path integrals and symmetry breaking for optimal control theory. *Journal of*  
308 *statistical mechanics: theory and experiment*, 2005(11):P11011.
- 309 Kappen, H. J., Gómez, V., and Opper, M. (2012). Optimal control as a graphical model inference  
310 problem. *Machine learning*, 87(2):159–182.
- 311 Kapturowski, S., Ostrovski, G., Quan, J., Munos, R., and Dabney, W. (2018). Recurrent experience  
312 replay in distributed reinforcement learning. In *International conference on learning representa-*  
313 *tions*.
- 314 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv e-prints*, pages  
315 arXiv–1412.
- 316 Leike, J., Martic, M., Krakovna, V., Ortega, P. A., Everitt, T., Lefrancq, A., Orseau, L., and Legg, S.  
317 (2017). AI safety gridworlds. *arXiv preprint arXiv:1711.09883*.
- 318 Markowitz, H. (1952). Portfolio selection. *Journal of Finance*, 1(7).
- 319 Mihatsch, O. and Neuneier, R. (2002). Risk-sensitive reinforcement learning. *Machine learning*,  
320 49(2):267–290.
- 321 Mikulik, V., Delétang, G., McGrath, T., Genewein, T., Martic, M., Legg, S., and Ortega, P. A. (2020).  
322 Meta-trained agents implement bayes-optimal agents. *arXiv preprint arXiv:2010.11223*.
- 323 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A.,  
324 Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep  
325 reinforcement learning. *nature*, 518(7540):529–533.
- 326 Nilim, A. and El Ghaoui, L. (2005). Robust control of Markov decision processes with uncertain  
327 transition matrices. *Operations Research*, 53(5):780–798.
- 328 Niv, Y., Edlund, J. A., Dayan, P., and O’Doherty, J. P. (2012). Neural prediction errors reveal  
329 a risk-sensitive reinforcement-learning process in the human brain. *Journal of Neuroscience*,  
330 32(2):551–562.
- 331 Ortega, D. A. and Braun, P. A. (2011). Information, utility and bounded rationality. In *International*  
332 *Conference on Artificial General Intelligence*, pages 269–274. Springer.
- 333 Ortega, P. A. and Braun, D. A. (2013). Thermodynamics as a theory of decision-making with  
334 information-processing costs. *Proceedings of the Royal Society A: Mathematical, Physical and*  
335 *Engineering Sciences*, 469(2153):20120683.
- 336 Ortega, P. A. and Lee, D. (2014). An adversarial interpretation of information-theoretic bounded  
337 rationality. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28.
- 338 Ortega, P. A., Wang, J. X., Rowland, M., Genewein, T., Kurth-Nelson, Z., Pascanu, R., Heess, N.,  
339 Veness, J., Pritzel, A., Sprechmann, P., et al. (2019). Meta-learning of sequential strategies. *arXiv*  
340 *preprint arXiv:1905.03030*.
- 341 Rescorla, R. A. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of  
342 reinforcement and nonreinforcement. *Current research and theory*, pages 64–99.

- 343 Robbins, H. and Monro, S. (1951). A stochastic approximation method. *The annals of mathematical*  
344 *statistics*, pages 400–407.
- 345 Russell, S., Dewey, D., and Tegmark, M. (2015). Research priorities for robust and beneficial artificial  
346 intelligence. *Ai Magazine*, 36(4):105–114.
- 347 Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. (2016). Meta-learning with  
348 memory-augmented neural networks. In *International conference on machine learning*, pages  
349 1842–1850. PMLR.
- 350 Sutton, R. S. and Barto, A. G. (1990). Time-derivative models of Pavlovian reinforcement. In  
351 *Learning and Computational Neuroscience: Foundations of Adaptive Networks*, pages 497–537.  
352 MIT Press.
- 353 Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press.
- 354 Tamar, A., Mannor, S., and Xu, H. (2014). Scaling up robust MDPs using function approximation. In  
355 *International Conference on Machine Learning*, pages 181–189. PMLR.
- 356 Theodorou, E., Buchli, J., and Schaal, S. (2010). A generalized path integral control approach to  
357 reinforcement learning. *The Journal of Machine Learning Research*, 11:3137–3181.
- 358 Tishby, N. and Polani, D. (2011). Information theory of decisions and actions. In *Perception-action*  
359 *cycle*, pages 601–636. Springer.
- 360 Todorov, E. (2007). Linearly-solvable Markov decision problems. In *Advances in neural information*  
361 *processing systems*, pages 1369–1376.
- 362 Toussaint, M. (2009). Robot trajectory optimization using approximate inference. In *Proceedings of*  
363 *the 26th annual international conference on machine learning*, pages 1049–1056.
- 364 Wang, J. X., Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., Blundell, C.,  
365 Kumaran, D., and Botvinick, M. (2016). Learning to reinforcement learn. *arXiv preprint*  
366 *arXiv:1611.05763*.
- 367 Ziebart, B. D., Maas, A. L., Bagnell, J. A., and Dey, A. K. (2008). Maximum entropy inverse  
368 reinforcement learning. In *AAAI*.

369 **Checklist**

- 370 1. For all authors...
- 371 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
372 contributions and scope? [Yes]
- 373 (b) Did you describe the limitations of your work? [Yes] see Section 5
- 374 (c) Did you discuss any potential negative societal impacts of your work? [Yes] see  
375 Section 5
- 376 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
377 them? [Yes]
- 378 2. If you are including theoretical results...
- 379 (a) Did you state the full set of assumptions of all theoretical results? [Yes] in Section 2.  
380 (b) Did you include complete proofs of all theoretical results? [Yes] in Section 2.
- 381 3. If you ran experiments...
- 382 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
383 imental results (either in the supplemental material or as a URL)? [No] The code is  
384 proprietary.
- 385 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
386 were chosen)? [Yes] see Section 4.
- 387 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
388 ments multiple times)? [Yes] see Section 4.
- 389 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
390 of GPUs, internal cluster, or cloud provider)? [Yes] see Section 4.
- 391 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 392 (a) If your work uses existing assets, did you cite the creators? [N/A]  
393 (b) Did you mention the license of the assets? [N/A]  
394 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
395
- 396 (d) Did you discuss whether and how consent was obtained from people whose data you're  
397 using/curating? [N/A]
- 398 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
399 information or offensive content? [N/A]
- 400 5. If you used crowdsourcing or conducted research with human subjects...
- 401 (a) Did you include the full text of instructions given to participants and screenshots, if  
402 applicable? [N/A]
- 403 (b) Did you describe any potential participant risks, with links to Institutional Review  
404 Board (IRB) approvals, if applicable? [N/A]
- 405 (c) Did you include the estimated hourly wage paid to participants and the total amount  
406 spent on participant compensation? [N/A]