

Adaptive stimulus selection for multi-alternative psychometric functions with lapses

Ji Hyun Bak^{1,2} and Jonathan W. Pillow³

¹School of Computational Sciences, Korea Institute for Advanced Study, Seoul, Korea; ²Department of Physics, Princeton University, NJ, USA;

³Department of Psychology and Princeton Neuroscience Institute, Princeton University, NJ, USA

(Dated: June 21, 2018)

Psychometric functions (PFs) quantify how external stimuli affect behavior and play an important role in building models of sensory and cognitive processes. Adaptive stimulus selection methods seek to select stimuli that are maximally informative about the PF given data observed so far in an experiment and thereby reduce the number of trials required to estimate the PF. Here we develop new adaptive stimulus selection methods for flexible PF models in tasks with two or more alternatives. We model the PF with a multinomial logistic regression mixture model that incorporates realistic aspects of psychophysical behavior, including lapses and multiple alternatives for the response. We propose an information-theoretic criterion for stimulus selection and develop computationally efficient methods for inference and stimulus selection based on semi-adaptive Markov Chain Monte Carlo (MCMC) sampling. We apply these methods to data from macaque monkeys performing a multi-alternative motion discrimination task, and show in simulated experiments that our method can achieve a substantial speed-up over random designs. These advances will reduce the data needed to build accurate models of multi-alternative PFs and can be extended to high-dimensional PFs that would be infeasible to characterize with standard methods.

Keywords: adaptive stimulus selection, sequential optimal design, Bayesian adaptive design, psychometric function, closed-loop experiments

1 Introduction

2 Understanding the factors governing psychophysical behavior is
3 a central problem in neuroscience and psychology. Although ac-
4 curate quantification of the behavior is an important goal in it-
5 self, psychophysics provides an important tool for interrogating
6 the mechanisms governing sensory and cognitive processing in the
7 brain. As new technologies allow direct manipulations of neural
8 activity in the brain, there is a growing need for methods that can
9 characterize rapid changes in psychophysical behavior.

10 In a typical psychophysical experiment, an observer is trained to
11 report judgements about a sensory stimulus by selecting a response
12 from among two or more alternatives. The observer is assumed to
13 have an internal probabilistic rule governing these decisions; this
14 probabilistic map from stimulus to response is called the observer's

15 psychometric function. Because the psychometric function is not
16 directly observable, it must be inferred from multiple observations
17 of stimulus-response pairs. However, such experiments are costly
18 due to the large numbers of trials typically required to obtain good
19 estimates of psychometric functions. Therefore, a problem of ma-
20 jor practical importance is to develop efficient experimental de-
21 signs that can minimize the amount of data required to accurately
22 infer an observer's psychometric function.

Bayesian adaptive stimulus selection. A powerful ap-
23 proach for improving the efficiency of psychophysical experiments
24 is to design the data collection process so that the stimulus is adap-
25 tively selected on each trial by maximizing a suitably defined ob-
26 jective function (MacKay, 1992). Such methods are known by a
27 variety of names, including “active learning”, “adaptive or sequen-
28 tial optimal experimental design”, and “closed-loop experiments.”
29

Bayesian approaches to adaptive stimulus selection define optimality of a stimulus in terms of its expected ability to improve the posterior distribution over the psychometric function, e.g., by reducing its variance or entropy. The three key ingredients of a Bayesian adaptive stimulus selection method are (Chaloner & Verdinelli, 1995; Pillow & Park, 2016):

- **model** - parametrizes the psychometric function of interest;
- **prior** - captures initial beliefs about model parameters;
- **utility function** - quantifies the usefulness of a hypothetical stimulus-response pair for improving the posterior.

Sequential algorithms for adaptive Bayesian experiments rely on repeated application of three basic steps: (i) data collection (stimulus presentation and response measurement); (ii) inference (posterior updating using data from the most recent trial); and (iii) selection of an optimal stimulus for the next trial by maximizing expected utility (see Fig. 1A). The inference step involves updating the posterior distribution over the model parameters according to Bayes rule with data from the most recent trial. Stimulus selection involves calculating the expected utility (i.e., the expected improvement in the posterior) for a set of candidate stimuli, averaging over the responses that might be elicited for each stimulus, and selecting the stimulus for which the expected utility is highest. Example utility functions include the negative trace of the posterior covariance (corresponding to the sum of the posterior variances for each parameter) and the mutual information or information gain (which corresponds to minimizing the entropy of the posterior).

Methods for Bayesian adaptive stimulus selection have been developed over several decades in a variety of different disciplines. If we focus on the specific application of estimating psychometric functions, the field goes back to the QUEST (A. B. Watson & Pelli, 1983) and ZEST (King-Smith, Grigsby, Vingrys, Benes, & Supowit, 1994) algorithms, which were focused on the estimation of discrimination thresholds, and to the simple case of 1-dimension stimulus and binary responses (Treutwein, 1995). The Ψ method (Kontsevich & Tyler, 1999) used Bayesian inference for estimat-

ing both threshold and slope of a psychometric function, which were extended to two-dimensional stimuli by Kujala and Lukka (2006). Further development of the method allowed for adaptive estimation of more complex psychometric functions, where the parameters were no longer limited to a threshold and a slope (Barthelmé & Mamassian, 2008; Kujala & Lukka, 2006; Lesmes, Lu, Baek, & Albright, 2010; Prins, 2013); and possibly related to each other (Vul, Bergsma, & MacLeod, 2010). Models with multi-dimensional stimuli were also considered (DiMattina, 2015; Kujala & Lukka, 2006; A. B. Watson, 2017).

Different models have been used to describe the psychometric function. Standard choices include the logistic regression model (Chaloner & Larntz, 1989; DiMattina, 2015; Zocchi & Atkinson, 1999), the Weibull distribution function (A. B. Watson & Pelli, 1983), and the cumulative function of Gaussian distribution (Kontsevich & Tyler, 1999). More recent works also considered Gaussian Process models (Gardner, Song, Weinberger, Barbour, & Cunningham, 2015). Most of the previous works, however, were limited to the case of binary responses.

In parallel, the development of Bayesian methods for inferring psychometric functions (Kuss, Jäkel, & Wichmann, 2005; Prins, 2012; Wichmann & Hill, 2001a, 2001b) have enlarged the space of statistical models that could be employed to describe psychophysical phenomena based on (often limited) data. A variety of recent advances also arose in sensory neuroscience or neurophysiology, driven by the development of efficient inference techniques for neural encoding models (Lewi, Butera, & Paninski, 2009; M. Park, Horwitz, & Pillow, 2011) or model comparison and discrimination methods (Cavagnaro, Myung, Pitt, & Kujala, 2010; DiMattina & Zhang, 2011; Kim, Pitt, Lu, Steyvers, & Myung, 2014). These advances can in many cases be equally well applied to psychophysical experiments.

Our contributions. In this paper, we develop methods for adaptive stimulus selection in psychophysical experiments that are applicable to realistic models of human and animal psychophysical behavior. First, we describe a psychophysical model that incor-

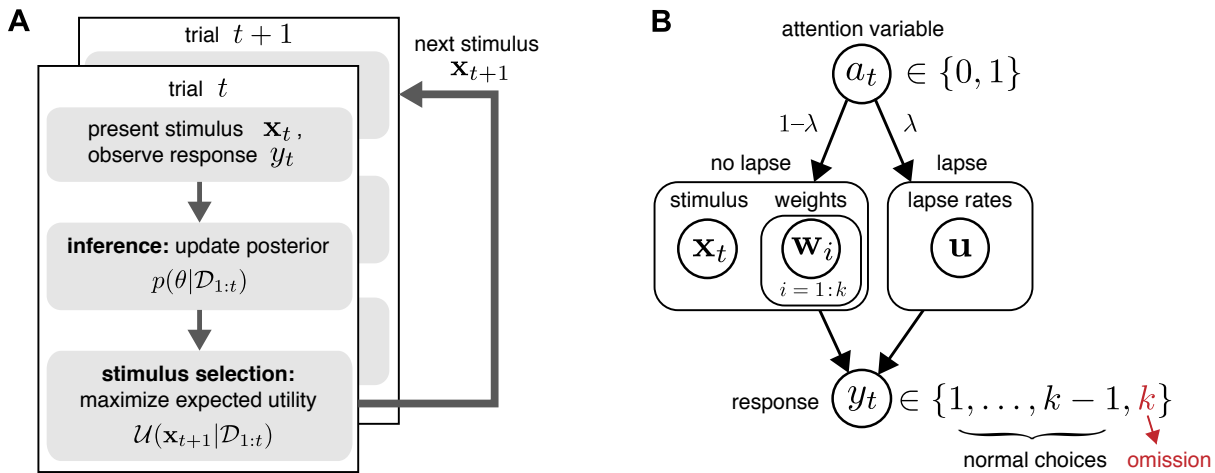


Figure 1: **(A)** Schematic of Bayesian adaptive stimulus selection. On each trial: (i) a stimulus is presented and response is observed; (ii) the posterior over the parameters θ is updated using all data collected so far in the experiment \mathcal{D}_t ; and (iii) the stimulus that maximizes the expected utility (in our case, information gain) is selected for the next trial. **(B)** A graphical model illustrating a hierarchical psychophysical observer model that incorporates lapses as well as the possibility of omissions. On each trial, a latent attention or lapse variable a_t is drawn from a Bernoulli distribution with parameter λ , to determine whether the observer attends to the stimulus \mathbf{x}_t on that trial or lapses. With probability $1 - \lambda$, the observer attends to the stimulus ($a_t = 0$), and the response y_t is drawn from a multinomial logistic regression model, where the probability of choosing option i is proportional to $\exp(\mathbf{w}_i^\top \mathbf{x}_t)$. With probability λ , the observer lapses ($a_t = 1$) and selects a choice from a (stimulus-independent) response distribution governed by parameter vector \mathbf{u} . So-called “omission” trials, in which the observer does not select one of the valid response options, are modeled with an additional response category $y_t = k$.

101 porates multiple response alternatives and “lapsing”, in which the
 102 observer makes a response that does not depend on the stimulus.
 103 This model can also incorporate “omission” trials, where the ob-
 104 server does not make a valid response (e.g., breaking fixation be-
 105 fore the go cue), by considering them as an additional response
 106 category. Second, we describe efficient methods for updating the
 107 posterior over the model parameters after every trial. Third, we in-
 108 troduce two algorithms for adaptive stimulus selection, one based
 109 on a Gaussian approximation to the posterior and a second based
 110 on Markov Chain Monte Carlo (MCMC) sampling. We apply these
 111 algorithms to simulated data and to real data analyzed with simu-
 112 lated closed-loop experiments, and show that they can substantially
 113 reduce in the number of trials required to estimate multi-alternative
 114 psychophysical functions.

115 Psychophysical observer model

116 Here we describe a flexible model of psychometric functions (PFs)
 117 based on the multinomial logistic (MNL) response model (Glonek
 118 & McCullagh, 1995). We show how omission trials can be nat-
 119 urally incorporated into a model as one of the multiple responses
 120 alternatives. We then develop a hierarchical extension of the model
 121 that incorporates lapses (see Fig. 1B).

122 **Multinomial logistic response model.** We consider the set-
 123 ting where the observer is presented with a stimulus $\mathbf{x} \in \mathbb{R}^d$ and
 124 selects a response $y \in \{1, \dots, k\}$ from one of k discrete choices on
 125 each trial. We will assume the stimulus is represented internally
 126 by some (possibly non-linear) feature vector $\phi(\mathbf{x})$, which we will
 127 write simply as ϕ for notational simplicity.

128 In the multinomial logistic model, the probability p_i of each pos-
 129 sible outcome $i \in \{1, \dots, k\}$ is determined by the dot product

between the feature ϕ and a vector of weights \mathbf{w}_i according to:

$$p_i = \frac{\exp(\mathbf{w}_i^\top \phi)}{\sum_{j=1}^k \exp(\mathbf{w}_j^\top \phi)}, \quad (1)$$

where the denominator ensures that these probabilities sum to 1, $\sum_{i=1}^k p_i = 1$. The function from stimulus to a probability vector over choices, $\mathbf{x} \mapsto (p_1, \dots, p_k)$, is the psychometric function, and the set of weights $\{\mathbf{w}_i\}_{i=1}^k$ are its parameters. Note that the model is over-parameterized when written this way, since the requirement that probabilities sum to 1 removes one degree of freedom from the probability vector. Thus, we can without loss of generality fix one of the weight vectors to zero, for example $\mathbf{w}_k = \mathbf{0}$, so that the denominator in (eq. 1) becomes $z = 1 + \sum_{j=1}^k \exp(\mathbf{w}_j^\top \phi)$ and $p_k = 1/z$.

We consider the feature vector ϕ to be a known function of the stimulus \mathbf{x} , even when the dependence is not written explicitly. For example, we can consider a simple form of feature embedding, $\phi(\mathbf{x}) = [1, \mathbf{x}^\top]^\top$, corresponding to a linear function of the stimulus plus an offset. In this case, the weights for the i 'th choice would correspond to $\mathbf{w}_i = [b_i, \mathbf{a}_i^\top]^\top$, where b_i is the offset or bias for the i 'th choice, and \mathbf{a}_i are the linear weights governing sensitivity to \mathbf{x} . The resulting choice probability has the familiar form, $p_i \propto \exp(b_i + \mathbf{a}_i^\top \mathbf{x})$. Nonlinear stimulus dependencies can be incorporated by including nonlinear functions of \mathbf{x} in the feature vector $\phi(\mathbf{x})$ (Knoblauch & Maloney, 2008; Murray, 2011; Neri & Heeger, 2002). Dependencies on the trial history, such as the previous stimulus or reward, may also be included as additional features in ϕ (see for example Bak, Choi, Akrami, Witten, and Pillow (2016)).

It is useful to always work with a normalized stimulus space, in which the mean of each stimulus component x_α over the stimulus space is $\langle x_\alpha \rangle = 0$, and the standard deviation $\text{std}(x_\alpha) = 1$. This normalization ensures that the values of the weight parameters are defined in more interpretable ways. The zero-mean condition ensures that the bias b is the expectation value of log probability over all possible stimuli. The unit-variance condition means that the effect of moving a certain distance along one dimension of the

weight space is comparable to the moving the same distance in another dimension, again averaged over all possible stimuli. In other words, we are justified to use the same unit along all dimensions of the weight space.

Including omission trials. Even in binary tasks with only two possible choices per trial, there is often an implicit third choice, which is to make no response, make an illegal response, or to interrupt the trial at some point before the allowed response period. For example, animals are often required to maintain an eye position or a nose poke, or wait for a “go” cue before reporting a choice. Trials on which the animal fails to obey these instructions are commonly referred to as “omissions” or “violations”, and are typically discarded from analysis. However, failure to take these trials into account may bias the estimates of the PF if they are more common for some stimuli than others (see Fig. 2B).

The multinomial response model provides a natural framework for incorporating omission trials because it accommodates an arbitrary number of response categories. Thus we can model omissions explicitly as one of the possible choices the observer can choose from, or as the $(k + 1)$ 'st response category in addition to the k valid responses. One could even consider different kinds of omissions separately, e.g., allowing choice $k+1$ to reflect fixation period violations and choice $k + 2$ to reflect failure to report a choice during the response window. Henceforth, we will let k reflect the total number of choices, including omission, as illustrated in Fig. 1B.

This formulation can also be useful for the rated Yes/No task in human psychophysics, where a “Not Sure” response is explicitly presented (C. S. Watson, Kellogg, Kawanishi, & Lucas, 1973). Although such model was considered for adaptive stimulus selection (Lesmes et al., 2015), the third alternative was not handled as a fully independent choice, as the goal was only to estimate the two detection thresholds separately: one for a strict Yes, another for a collapsed response of either Yes or Not Sure. Our model treats each of the multiple alternatives equivalently.

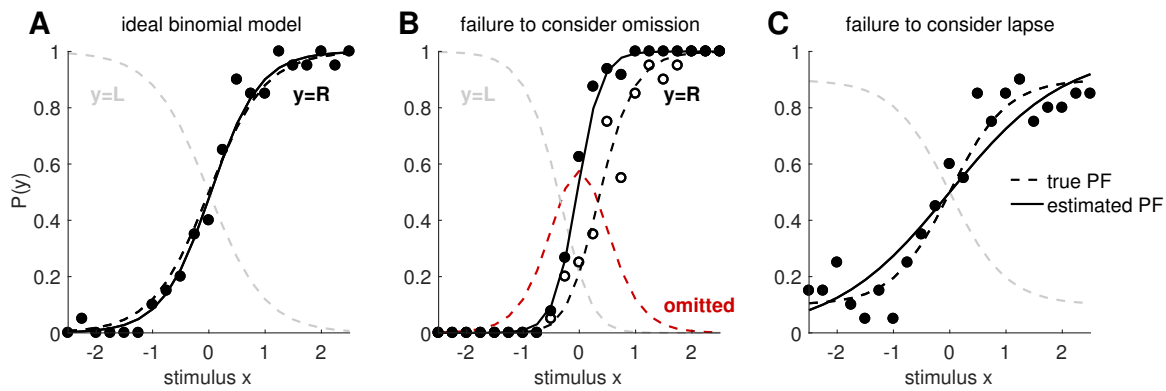


Figure 2: Effects of omission and lapse. Here we illustrate the undesirable effects of failing to take into account of omission and lapse. **(A)** If the PF follows an ideal binomial logistic model, it can be estimated very well from data. The black dashed line shows the true PF for one of the two responses (say $y = R$), and the gray dashed line shows the true PF for the other response (say $y = L$), such that the two dashed curves always add up to 1. The black dots indicate the mean probability to observe this response $y = R$ at each stimulus point x . We drew 20 observations per stimulus point, at each of the 21 stimulus points along the 1-dimensional axis. The resulting estimate for $P(y = 1)$ is shown in the solid black line. The inference method is not important for the current purpose, but we used the MAP estimate, discussed in a later section. **(B)** Now suppose that some trials fell into the implicit third choice which is the omission (red dashed line shows omission probability). The observed probability of $y = R$ at each stimulus point (open black circles) follows the true PF (black dashed line). But if the omitted trials are systematically excluded from analysis, as in common practice, the estimated PF (solid black line) reflects a biased set of observations (filled black circles), and fail to recover the true PF. **(C)** When there is a finite lapse rate (we used a total lapse of $\lambda = 0.2$, uniformly distributed to the two outcomes), the true PF (dashed black line) asymptotes to a finite offset from 0 or 1. If the resulting observations (black dots) are fitted to a plain binomial model without lapse, the slope of the estimated PF (solid black line) is systematically biased.

199 **Modeling lapse with a mixture model.** Another important
 200 feature of real psychophysical observers is the tendency to occa-
 201 sionally make errors that are independent of the stimulus. Such
 202 choices, commonly known as “lapses” or “button press errors”,
 203 may reflect lapses in concentration or memory of the response
 204 mapping (Kuss et al., 2005; Wichmann & Hill, 2001a). Lapses
 205 are most easily identified by errors on “easy” trials, that is, trials
 206 that should be performed perfectly if the observer were paying at-
 207 tention.

208 Although lapse rates can be negligible in highly trained obser-
 209 vers (Carandini & Churchland, 2013), they can be substantially
 210 greater than zero in settings involving non-primates or complicated
 211 psychophysical tasks. Lapses affect the psychometric function by
 212 causing it to saturate above 0 and below 1, so that “perfect” per-
 213 formance is never achieved even for the easiest trials. Failure to
 214 incorporate lapses into the PF model may therefore bias estimates
 215 of sensitivity, as quantified by PF slope or threshold (illustrated

in Fig. 2C; also see Wichmann and Hill (2001a, 2001b) or Prins
 (2012)).

To model lapses, we use a mixture model that treats the obser-
 ver’s choice on each trial as coming from one of two probability
 distributions: a stimulus-dependent distribution (governed by the
 multinomial logistic model) and stimulus-independent distribution
 (reflecting a fixed probability of choosing any option when laps-
 ing). Simpler versions of such mixture model have been proposed
 previously (Kuss et al., 2005).

Fig. 1B shows a schematic of the resulting model. On each trial,
 a Bernoulli random variable $a \sim \text{Ber}(\lambda)$ governs whether the ob-
 server lapses: with probability λ and the observer lapses (i.e., ig-
 nores the stimulus), and with probability $1 - \lambda$, and the observer at-
 tends to the stimulus. If the observer lapses ($a = 1$), the response is
 drawn according to fixed probability distribution (c_1, \dots, c_k) gov-
 erning the probability of selecting options 1 to k , where $\sum c_i = 1$.
 If the observer does not lapse ($a = 0$), the observer selects a re-

233 sponse according to the multinomial logistic model. Under this
234 model, the conditional probability of choosing option i given the
235 stimulus can be written:

$$236 \quad p_i = (1 - \lambda)q_i + \lambda c_i, \quad q_i = \frac{\exp(\mathbf{w}_i^\top \phi)}{\sum_j \exp(\mathbf{w}_j^\top \phi)} \quad (2)$$

237 where q_i is the lapse-free probability under the classical
238 MNL model (eq. 1).

239 It is convenient to re-parameterize this model so that λc_i , the
240 conditional probability of choosing the i -th option due to a lapse,
241 is written

$$242 \quad \lambda c_i = \frac{\exp(u_i)}{1 + \sum_j \exp(u_j)}, \quad (3)$$

243 where each auxiliary lapse parameter u_i is proportional to the log
244 probability of choosing option i due to lapse. The lapse-conditional
245 probabilities of each choice, c_i , and the total lapse probability, λ ,
246 are respectively

$$247 \quad c_i = \frac{\exp(u_i)}{\sum_j \exp(u_j)}, \quad \lambda = \sum_i \frac{\exp(u_i)}{1 + \sum_j \exp(u_j)}. \quad (4)$$

248 Because each u_i lives on the entire real line, fitting can be carried
249 out with unconstrained optimization methods, although adding rea-
250 sonable constraints may improve performance in some cases. The
251 full parameter vector of the resulting model is $\boldsymbol{\theta} = [\mathbf{w}^\top, \mathbf{u}^\top]^\top$,
252 which includes k additional lapse parameters $\mathbf{u} = \{u_1, \dots, u_k\}$.
253 Note that in some cases it might be desirable to assume lapse
254 choices obey a uniform distribution, where the probability of each
255 option is $c_i = 1/k$. For this simplified “uniform-lapse” model we
256 need only a single lapse parameter u . Note that we have unified
257 the parameterizations of the “lapse rate” (deviation of the upper
258 asymptote of the PF from 1; in this case $\lambda - \lambda c_i$) and the “guess
259 rate” (deviation of the lower asymptote from 0; in this case λc_i),
260 which are often modeled separately in previous works with two-
261 alternative responses (Schütt, Harmeling, Macke, & Wichmann,
262 2016; Wichmann & Hill, 2001a, 2001b). Here they are written in
263 terms of a single family of parameters $\{u_i\}$, and extended naturally
264 to multi-alternative responses.

265 Our model provides a general and practical parametrization of
266 psychometric functions with lapses. Although previous work has

267 considered the problem of modeling lapses in psychophysical data,
268 much of it assumed a uniform-lapse model, where all options are
269 equally likely during lapses. Earlier approaches have often
270 assumed either that the lapse probability was known a priori (Kontse-
271 vich & Tyler, 1999), or was fit by a grid search over a small set of
272 candidate values (Wichmann & Hill, 2001a). Here we instead in-
273 fer individual lapse probabilities for each response option, similar
274 to recent approaches described in Kuss et al. (2005); Prins (2012,
275 2013); Schütt et al. (2016). Importantly, our method infers the full
276 parameter $\boldsymbol{\theta}$ that includes both the weight and the lapse parameters,
277 rather than treating the lapse separately. In particular, our parame-
278 terization (eq. 3) has the advantage that there is no need to constrain
279 the support of the lapse parameters u_i . These parameters’ relation-
280 ship to lapse probabilities c_i takes the same (“softmax”) functional
281 form as the multinomial logistic model, placing both sets of param-
282 eters on an equal footing.

283 Before closing this section, we would like to reflect briefly on
284 the key differences between omissions and lapses. First, although
285 omissions and lapses both reflect errors in decision making, omis-
286 sions are defined as invalid responses and are thus easily identi-
287 fiable from the data; lapses, on the other hand, are indistinguish-
288 able from normal responses, and are identifiable only from the fact
289 that the psychometric function does not saturate at 0 or 1. Second,
290 modeling omissions as a response category under the multinomial
291 logistic model means that the probability of omission is stimulus-
292 dependent (e.g., more likely to arise on trials with high difficulty, or
293 generally when the evidence for other options is low). Even if the
294 omissions are not stimulus-dependent, and governed entirely by a
295 “bias” parameter, the probability of omission will still be higher
296 when the evidence for other choices is low, or lower when the ev-
297 idence for other choices is high. Omissions that arise in a purely
298 stimulus-independent fashion, on the other hand, will be modeled
299 as arising from the lapse parameter associated with the omission
300 response category. Omissions can thus arise in two ways under the
301 model: as categories selected under the multinomial model or as
302 lapses arising independent of the stimulus and other covariates.

Posterior inference

Bayesian methods for adaptive stimulus selection require the posterior distribution over model parameters given the data observed so far in an experiment. The posterior distribution results from the combination of two ingredients: a prior distribution $p(\boldsymbol{\theta})$, which captures prior uncertainty about the model parameters $\boldsymbol{\theta}$, and a likelihood function $p(\{y_s\}|\{\mathbf{x}_s\}, \boldsymbol{\theta})$, which captures information about the parameters from the data $\{(\mathbf{x}_s, y_s)\}$, $s = 1, \dots, t$, consisting of stimulus-response pairs observed up to the current time bin t .

Unfortunately, the posterior distribution for our model has no analytic form. We therefore describe two methods for approximate posterior inference: one relying on a Gaussian approximation to the posterior, known as the Laplace approximation, and a second one based on MCMC sampling.

Prior. The prior distribution specifies our beliefs about model parameters before we have collected any data, and serves to regularize estimates obtained from small amounts of data, e.g., by shrinking estimated weights toward zero. Typically we want the prior to be weak enough that the likelihood dominates the posterior for reasonable-sized datasets. However, the choice of prior is especially important in adaptive stimulus selection settings because it determines the effective volume of the search space (M. Park & Pillow, 2012; M. Park, Weller, Horwitz, & Pillow, 2014). For example, if the weights are known to exhibit smoothness, then a correlated or smoothness-inducing prior can improve the performance of adaptive stimulus selection because the effective size (or entropy) of the parameter space is much smaller than under an independent prior (M. Park & Pillow, 2012).

In this study, we use a generic independent, zero-mean Gaussian prior over the weight vectors

$$p(\mathbf{w}_i) = \mathcal{N}(\mathbf{0}, \sigma^2 I), \quad (5)$$

for all $i \in (1, \dots, k)$, with a fixed standard deviation σ . This choice of prior is appropriate when the regressors $\{\mathbf{x}\}$ are standardized,

since any single weight can take values that allow for a range of psychometric function shapes along that axis, from flat ($w = 0$) to steeply decreasing ($w = -2\sigma$) or increasing ($w = +2\sigma$). We used $\sigma = 3$ in the simulated experiments in Results. For the lapse parameters $\{u_i\}$, we used a uniform prior over the range $[\log(0.001), 0]$ with the natural log, so that each lapse probability λ_{c_i} is bounded between 0.001 and $1/2$. We set the lower range constraint below $1/N$, where $N = 100$ is the number of observed trials in our simulations, since we cannot reasonably infer lapse probabilities with precision finer than $1/N$. The upper range constraint gives maximal lapse probabilities of $1/(k+1)$ if all u_i take on the maximal value of 0. Note that our prior is uniform with respect to the rescaled lapse parameters $\{u_i\}$, rather than to the actual lapse rates; projected to the space of the lapse probabilities, given the bounds, the prior increases towards smaller lapse. For a comprehensive study of the effect of different priors on lapse, see for example Schütt et al. (2016).

Psychometric function likelihood. The likelihood is the conditional probability of the data as a function of the model parameters. Although we have thus far considered the response variable y to be a scalar taking values in the set $\{1, \dots, k\}$, it is more convenient to use a “one-hot” or “1-of- k ” representation, in which the response variable \mathbf{y} for each trial is a length- k vector with one 1 and $(k-1)$ zeros; the position of the 1 in this vector indicates the category selected. For example, in a task with four possible options per trial, a response vector $\mathbf{y} = [0 \ 0 \ 1 \ 0]$ indicates a trial on which the observer selected the third option.

With this parametrization, the log-likelihood function for a single trial can be written

$$\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \sum_i y_i \log p_i(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{y}^\top \log \mathbf{p}(\mathbf{x}, \boldsymbol{\theta}), \quad (6)$$

where $p_i(\mathbf{x}, \boldsymbol{\theta})$ denotes the probability $p(y_i = 1|\mathbf{x}, \boldsymbol{\theta})$ under the model (eq. 1), and $\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}) \equiv [p_1(\mathbf{x}, \boldsymbol{\theta}), \dots, p_k(\mathbf{x}, \boldsymbol{\theta})]^\top$ denotes the vector of probabilities for a single trial.

In the classical (lapse-free) multinomial logistic model, where $\boldsymbol{\theta} = \{\mathbf{w}_i\}$, the log likelihood is a concave function of $\boldsymbol{\theta}$, which

372 guarantees that numerical optimization of the log-likelihood will
 373 find a global optimum. With a finite lapse rate, however, the log
 374 likelihood is no longer concave. (See [Appendix A](#)).

375 **Posterior distribution.** The log-posterior can be written as the
 376 sum of log-prior and log-likelihood summed over trials, plus a con-
 377 stant:

$$378 \log p(\boldsymbol{\theta}|\mathcal{D}_t) = \log p(\boldsymbol{\theta}) + \sum_{s=1}^t \log p(y_s|\mathbf{x}_s, \boldsymbol{\theta}) + c, \quad (7)$$

379 where $\mathcal{D}_t \equiv \{\mathbf{x}_s, y_s\}_{s=1}^t$ denotes the accumulated data up to trial
 380 t and $c = -\log(\int p(\boldsymbol{\theta}) \prod_s p(y_s|\mathbf{x}_s) d\boldsymbol{\theta})$ is a normalization con-
 381 stant that does not depend on the parameters $\boldsymbol{\theta}$. Because this con-
 382 stant has no tractable analytic form, we rely on two alternate meth-
 383 ods for obtaining a normalized posterior distribution.

384 **Inference via Laplace approximation.** The Laplace approx-
 385 imation is a well-known Gaussian approximation to the posterior
 386 distribution, which can be derived from a second-order Taylor se-
 387 ries approximation to the log-posterior around its mode ([Bishop,](#)
 388 [2006](#)).

389 Computing the Laplace approximation involves a two-step pro-
 390 cedure. The first step is to perform a numerical optimization of
 391 $\log p(\boldsymbol{\theta}|\mathcal{D}_t)$ to find the posterior mode, or maximum a posteriori
 392 (MAP) estimate of $\boldsymbol{\theta}$. This vector, given by

$$393 \hat{\boldsymbol{\theta}}_t = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log p(\boldsymbol{\theta}) + \sum_{s=1}^t \log p(y_s|\mathbf{x}_s, \boldsymbol{\theta}), \quad (8)$$

394 provides the mean of the Laplace approximation. Because we can
 395 explicitly provide the gradient and Hessian of the log likelihood
 396 (see [Appendix A](#)) and log-prior, this optimization can be carried
 397 efficiently via Newton-Raphson or trust region methods.

398 The second step is to compute the second derivative (the Hes-
 399 sian matrix) of the log-posterior at the mode, which provides the
 400 inverse covariance of the Gaussian. This gives us a local Gaussian
 401 approximation of the posterior, centered at the posterior mode:

$$402 p(\boldsymbol{\theta}|\mathcal{D}_t) \approx \mathcal{N}(\hat{\boldsymbol{\theta}}_t, C_t), \quad (9)$$

403 where covariance $C_t = -H_t^{-1}$ is the inverse Hessian of the log
 404 posterior, $H_t(i, j) = \partial^2(\log p(\boldsymbol{\theta}|\mathcal{D}_t))/(\partial\theta_i\partial\theta_j)$, evaluated at $\hat{\boldsymbol{\theta}}_t$.

Note that when the log-posterior is concave (i.e., when the
 model does *not* contain lapse), numerical optimization is guaran-
 teed to find a global maximum of the posterior. Log-concavity
 also strengthens the rationale for using the Laplace approximation,
 since the true and approximate posterior are both log-concave den-
 sities centered on the true mode ([Paninski et al., 2010](#); [Pillow, Ah-
 madian, & Paninski, 2011](#)). When the model incorporates lapses,
 these guarantees no longer apply globally.

Inference via MCMC sampling. A second approach to in-
 ference is to generate samples from the posterior distribution over
 the parameters via Markov Chain Monte Carlo (MCMC) sampling.
 Sampling-based methods are typically more computationally in-
 tensive than the Laplace approximation, but may be warranted
 when the posterior is not provably log-concave (as is the case when
 lapse rates are non-zero) and therefore not well approximated by a
 single Gaussian.

The basic idea in MCMC sampling is to set up an easy-to-sample
 Markov Chain that has the posterior as its stationary distribution.
 Sampling from this chain produces a dependent sequence of pos-
 terior samples: $\{\boldsymbol{\theta}_m\} \sim p(\boldsymbol{\theta}|\mathcal{D}_t)$, which can be used to evaluate
 posterior expectations via Monte Carlo integrals:

$$426 \mathbb{E}[f(\boldsymbol{\theta})] \approx \frac{1}{M} \sum_{m=1}^M f(\boldsymbol{\theta}_m), \quad (10)$$

427 for any function $f(\boldsymbol{\theta})$. The mean of the posterior is obtained from
 428 setting $f(\boldsymbol{\theta}) = \boldsymbol{\theta}$, although for adaptive stimulus selection we will
 429 be interested in the full shape of the posterior.

The Metropolis-Hastings (MH) algorithm is perhaps the sim-
 plest and most widely-used MCMC sampling method ([Metropo-
 lis, Rosenbluth, Rosenbluth, Teller, & Teller, 1953](#)). It generates
 samples via a proposal distribution centered on the current sample
 (see [Appendix B](#)). The choice of proposal distribution is critical to
 the efficiency of the MH algorithm, since this governs the rate of
 “mixing”, or the the number of Markov Chain samples required to
 obtain independent samples from the posterior distribution ([Rosen-
 thal, 2011](#)). Faster mixing implies that fewer samples M are re-
 quired to obtain an accurate approximation to the posterior.

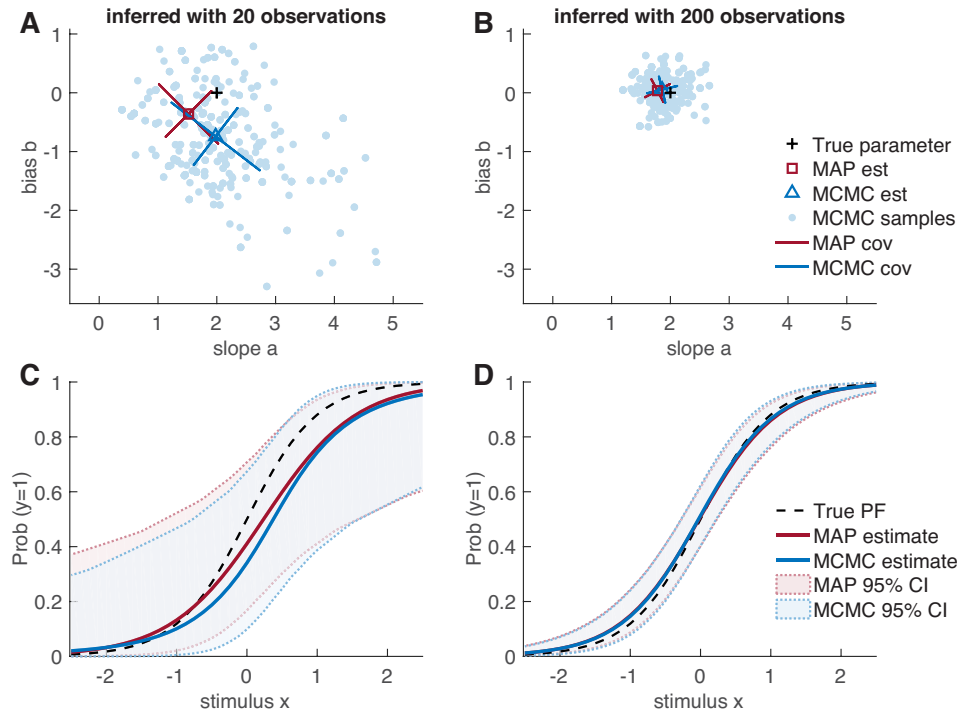


Figure 3: Inferring the psychometric function. Example of a psychometric problem, with a lapse-free binomial logistic model $f(v) = e^v / (1 + e^v)$. Given a 1D stimulus, a response were drawn from a “true” model $P(y = 1) = f(b + ax)$ with two parameters, slope $a = 2$ and bias $b = 0$. **(A-B)** Viewing on the parameter space, the posterior distributions become sharper (and closer to the true parameter values) as the dataset size N increases. Shown at a small **(A)** $N = 20$ and a large **(B)** $N = 200$. For the MAP estimate, the mode of the distribution is marked with a square, and the two standard deviations (“widths”) of its Gaussian approximation are shown with bars. For the MCMC sampling method, all $M = 500$ samples of the chain are shown in dots, the sample mean with a triangle, and the widths with the bars. The widths are the standard deviations along the principal directions of the sampled posterior (eigenvectors of the covariance matrix; not necessary aligned with the $a - b$ axes). **(C-D)** The accuracy of the estimated PF improves with the number of observations N , using either of the two posterior inference methods (MAP-based and sampling-based). Shown at a small **(C)** $N = 20$ and a large **(D)** $N = 200$. The two methods are highly consistent in this simple case, especially when N is large enough.

440 Here we propose a semi-adaptive MH algorithm, developed
 441 specifically for the current context of sequential learning. Our
 442 approach is based on an established observation that the optimal
 443 width of the proposal distribution should be proportional to the
 444 typical length scale of the distribution being sampled (Gelman,
 445 Roberts, & Gilks, 1996; Roberts, Gelman, & Gilks, 1997). Our al-
 446 gorithm is motivated by the adaptive Metropolis algorithm (Haario,
 447 Saksman, & Tamminen, 2001), where the proposal distribution is
 448 updated at each proposal within a single chain; here we do not
 449 adapt the proposal within chains, but rather after each trial. Specif-

ically, we set the covariance of a Gaussian proposal distribution to
 be proportional to the covariance of the samples from the previous
 trial, using the scaling factor of Haario et al. (2001). See Appendix
 B for details. The adaptive algorithm takes advantage of the fact
 that the posterior cannot change too much between trials, since it
 changes only by a single-trial likelihood term on each trial.

Adaptive stimulus selection methods

As data are collected during the experiment, the posterior distribution becomes narrower due to the fact that each trial carries some additional information about the model parameters (see Fig. 3). This narrowing of the posterior is directly related to information gain. A stimulus that produces no expected narrowing of the posterior is, by definition, uninformative about the parameters. On the other hand, a stimulus that (on average) produces a large change in the current posterior is an informative stimulus. Selecting informative stimuli will reduce the number of stimuli required to obtain a narrow posterior, which is the essence of adaptive stimulus selection methods. In this section, we introduce a precise measure of information gain between a stimulus and the model parameters, and propose an algorithm for selecting stimuli to maximize it.

Infomax criterion for stimulus selection. At each trial, we present a stimulus \mathbf{x} and observe the outcome \mathbf{y} . After t trials, the expected gain in information from a stimulus \mathbf{x} is equal to the mutual information between \mathbf{y} and the model parameters $\boldsymbol{\theta}$, given the data \mathcal{D}_t observed so far in the experiment. We denote this conditional mutual information:

$$I_t(\boldsymbol{\theta}; \mathbf{y}|\mathbf{x}) = \iint d\boldsymbol{\theta} d\mathbf{y} p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{x}, \mathcal{D}_t) \log \frac{p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{x}, \mathcal{D}_t)}{p(\boldsymbol{\theta}|\mathcal{D}_t)p(\mathbf{y}|\mathbf{x}, \mathcal{D}_t)}, \quad (11)$$

where $p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{x}, \mathcal{D}_t)$ is the joint distribution of $\boldsymbol{\theta}$ and \mathbf{y} given a stimulus \mathbf{x} and dataset \mathcal{D}_t , the term $p(\boldsymbol{\theta}|\mathcal{D}_t)$ is the current posterior distribution over the parameters from previous trials, and $p(\mathbf{y}|\mathbf{x}, \mathcal{D}_t) = \int d\boldsymbol{\theta} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D}_t)$ is known as the posterior-predictive distribution of \mathbf{y} given \mathbf{x} .

It is useful to note that the mutual information can equivalently be written in two other ways involving Shannon entropy. The first is given by:

$$I_t(\boldsymbol{\theta}; \mathbf{y}|\mathbf{x}) = H_t(\boldsymbol{\theta}) - H_t(\boldsymbol{\theta}|\mathbf{y}; \mathbf{x}) \quad (12)$$

where the first term is the entropy of the posterior at time t ,

$$H_t(\boldsymbol{\theta}) = - \int d\boldsymbol{\theta} p(\boldsymbol{\theta}|\mathcal{D}_t) \log p(\boldsymbol{\theta}|\mathcal{D}_t), \quad (13)$$

and the second is the conditional entropy of $\boldsymbol{\theta}$ given \mathbf{y} ,

$$\begin{aligned} H_t(\boldsymbol{\theta}|\mathbf{y}; \mathbf{x}) &= -\mathbb{E}_{\boldsymbol{\theta}, \mathbf{y}} \left[\log p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, \mathcal{D}_t) \right] \\ &= - \iint d\boldsymbol{\theta} d\mathbf{y} p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{x}, \mathcal{D}_t) \log p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, \mathcal{D}_t), \end{aligned} \quad (14)$$

which is the entropy of the updated posterior *after* having observed \mathbf{x} and \mathbf{y} , averaged over draws of \mathbf{y} from the posterior predictive distribution. Written this way, the mutual information can be seen as the expected reduction in posterior entropy from a new stimulus-response pair. Moreover, the first term, $H_t(\boldsymbol{\theta})$, is independent of the stimulus and response on the current trial, so infomax stimulus selection is equivalent to picking the stimulus that minimizes the expected posterior entropy $H_t(\boldsymbol{\theta}|\mathbf{y}; \mathbf{x})$.

A second equivalent expression for the mutual information, which will prove useful for our sampling-based method, is:

$$I_t(\boldsymbol{\theta}; \mathbf{y}|\mathbf{x}) = H_t(\mathbf{y}; \mathbf{x}) - H_t(\mathbf{y}|\boldsymbol{\theta}; \mathbf{x}), \quad (15)$$

which is the difference between the marginal entropy of the response distribution conditioned on \mathbf{x} ,

$$H_t(\mathbf{y}; \mathbf{x}) = - \int d\mathbf{y} p(\mathbf{y}|\mathbf{x}, \mathcal{D}_t) \log p(\mathbf{y}|\mathbf{x}, \mathcal{D}_t) \quad (16)$$

and the conditional entropy of the response \mathbf{y} given $\boldsymbol{\theta}$, conditioned on the stimulus:

$$H_t(\mathbf{y}|\boldsymbol{\theta}; \mathbf{x}) = - \iint d\mathbf{y} d\boldsymbol{\theta} p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{x}, \mathcal{D}_t) \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}). \quad (17)$$

This formulation shows the mutual information to be equal to the difference between the entropy of the marginal distribution of \mathbf{y} conditioned on \mathbf{x} (with $\boldsymbol{\theta}$ integrated out) and the average entropy of \mathbf{y} given \mathbf{x} and $\boldsymbol{\theta}$, averaged over the posterior distribution of $\boldsymbol{\theta}$. The dual expansion of the mutual information was also used by Kujala and Lukka (2006).

In a sequential setting where t is the latest trial and $t + 1$ is the upcoming one, the optimal stimulus is the information-maximizing (“infomax”) solution:

$$\mathbf{x}_{t+1} = \arg \max_{\mathbf{x}} I_t(\boldsymbol{\theta}; \mathbf{y}|\mathbf{x}). \quad (18)$$

Fig. 4 shows an example of a simulated experiment where the stimulus was selected adaptively following the infomax criterion. Note

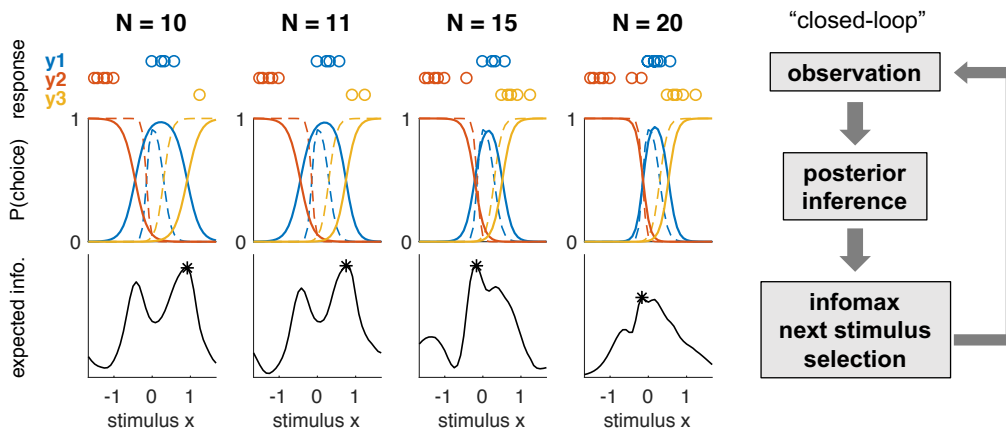


Figure 4: Example of infomax adaptive stimulus selection, simulated with a three-alternatives lapse-free model on 1D stimulus. The figure shows how given a small set of data (the stimulus-response pairs shown in top row), the PFs are estimated based on the accumulated data (middle row), and the next stimulus is chosen to maximize the expected information gain (bottom row). Each column shows the instance after the N observations in a single adaptive stimulus selection sequence, for $N = 10, 11, 15$ and 20 respectively. In the middle row, the estimated PFs (solid lines) quickly approach the true PFs (dashed lines) through the adaptive and optimal selection of stimuli. This example was generated using the Laplace approximation based algorithm, with an independent Gaussian prior over the weights with mean zero and standard deviation $\sigma = 10$.

525 that our algorithm takes a “greedy” approach of optimizing one
 526 trial at a time. For work on optimizing beyond the next trial, see
 527 for example [Kim, Pitt, Lu, and Myung \(2017\)](#).

528 Selecting the optimal stimulus thus requires maximizing the mu-
 529 tual information over the set of all possible stimuli $\{\mathbf{x}\}$. Since each
 530 evaluation of the mutual information involves a high-dimensional
 531 integral over parameter space and response space, this is a highly
 532 computationally demanding task. In the next sections, we present
 533 two algorithms for efficient infomax stimulus selection based on
 534 each of the two approximate inference methods described previ-
 535 ously.

536 **Infomax with Laplace approximation.** Calculation of the
 537 mutual information is greatly simplified by a Gaussian approxima-
 538 tion of the posterior. The entropy of a Gaussian distribution with
 539 covariance C is equal to $\frac{1}{2} \log |C|$ up to a constant factor. If we ex-
 540 pand the mutual information as in (eq. 12), and recall that we need
 541 only minimize the expected posterior entropy after observing the
 542 response, the optimal stimulus for time-step $t + 1$ is given by:

$$543 \mathbf{x}_{t+1}^* = \operatorname{argmin}_{\mathbf{x}} \int d\mathbf{y} p(\mathbf{y}|\mathbf{x}, \mathcal{D}_t) \log |\tilde{C}(\mathbf{x}, \mathbf{y})|, \quad (19)$$

544 where $\tilde{C}(\mathbf{x}, \mathbf{y})$ is the covariance of the updated (Gaussian) poste-
 545 rior after observing stimulus-response pair (\mathbf{x}, \mathbf{y}) . To evaluate the
 546 updated covariance $\tilde{C}(\mathbf{x}, \mathbf{y})$ under the Laplace approximation, we
 547 would need to numerically optimize the posterior for θ for each
 548 possible response \mathbf{y} , for any candidate stimulus \mathbf{x} , which would be
 549 computationally infeasible. We therefore use a fast approximate
 550 method for obtaining a closed-form update for $\tilde{C}(\mathbf{x}, \mathbf{y})$ from the
 551 current posterior covariance C_t , following an approach developed
 552 in [Lewi et al. \(2009\)](#). See [Appendix C](#) for details. Note that this
 553 approximate sequential update is only used for calculating the ex-
 554 pected utility of each candidate stimulus by approximating the pos-
 555 terior distribution at the next trial. For obtaining the MAP estimate
 556 of the current model parameter, θ_t , numerical optimization needs
 557 to be performed using the full accumulated data \mathcal{D}_t each time.

558 Once we have $\log |\tilde{C}(\mathbf{x}, \mathbf{y})|$ for each given stimulus-observation
 559 pair, we numerically sum this over a set of discrete counts \mathbf{y} that
 560 are likely under the posterior-predictive distribution. This is done

561 in two steps, by separating the integral in (eq. 19) as:

$$562 \int d\mathbf{y} p(\mathbf{y}|\mathbf{x}, \mathcal{D}_t) \log |\tilde{C}(\mathbf{x}, \mathbf{y})|$$

$$563 = \int d\boldsymbol{\theta}_t p(\boldsymbol{\theta}_t|\mathcal{D}_t) \int d\mathbf{y} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_t) \log |\tilde{C}(\mathbf{x}, \mathbf{y})|. \quad (20)$$

564 Note that the outer integral is over the current posterior $p(\boldsymbol{\theta}_t|\mathcal{D}_t) \approx$
 565 $\mathcal{N}(\hat{\boldsymbol{\theta}}_t, C_t)$, which is to be distinguished from the future posterior
 566 $p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, \mathcal{D}_t) \approx \mathcal{N}(\tilde{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}), \tilde{C}(\mathbf{x}, \mathbf{y}))$ whose entropy we are trying
 567 to minimize. Whereas the inner integral is simply a weighted sum
 568 over the set of outcomes \mathbf{y} , the outer integral over the parameter $\boldsymbol{\theta}$
 569 is in general challenging, especially when the parameter space is
 570 high-dimensional. In the case of the standard multinomial logistic
 571 model that does not include lapse, we can exploit the linear struc-
 572 ture of model to reduce this to a lower-dimensional integral over
 573 the space of the linear predictor, which we evaluate numerically
 574 using Gauss-Hermite quadrature (Heiss & Winschel, 2008). (This
 575 integral is 1D for classic logistic regression, and $(k-1)$ -dimensional
 576 for multinomial logistic regression with k classes; see Appendix
 577 C for details.) When the model incorporates lapses, the full pa-
 578 rameter vector $\boldsymbol{\theta} = [\mathbf{w}^\top, \mathbf{u}^\top]^\top$ includes the lapse parameters in
 579 addition to the weights \mathbf{w} . In this case, our method with Laplace
 580 approximation may suffer from reduced accuracy due to the fact
 581 that the posterior may be less closely approximated by a Gaussian.

582 In order to exploit the convenient structure of reduced integral
 583 over the weight space, we choose to maximize the *partial* infor-
 584 mation between the observation and the psychophysical weights,
 585 $I(\mathbf{w}; \mathbf{y}|\mathbf{x})$, instead of the full information $I(\boldsymbol{\theta}; \mathbf{y}|\mathbf{x})$. This is
 586 a reasonable approximation in many cases where the stimulus-
 587 dependent behavior is the primary focus of the psychometric ex-
 588 periment (also see Prins (2013) for a similar approach). How-
 589 ever, we note that this is the only piece in this work where we
 590 treat the weights separately from the lapse parameters; posterior
 591 inference is still performed for the full parameter $\boldsymbol{\theta}$. Thus for
 592 Laplace-based infomax exclusively, the partial covariance $C_{\mathbf{w}\mathbf{w}} =$
 593 $-(\partial^2(\log \mathcal{P})/\partial \mathbf{w}^2)^{-1}$ is used in place of the full covariance
 594 $C = -(\partial^2(\log \mathcal{P})/\partial \boldsymbol{\theta}^2)^{-1}$, where $\mathcal{P}(\boldsymbol{\theta})$ is the posterior distri-
 595 bution over the full parameter space. Because the positive semi-

596 definiteness of the partial covariance is still not guaranteed, it needs
 597 to be approximated to the nearest symmetric positive semi-definite
 598 matrix when necessary (Higham, 1988). We can show, however,
 599 that the partial covariance is asymptotically positive semi-definite
 600 in the small-lapse limit (Appendix A).
 601

602 **Infomax with MCMC.** Sampling-based inference provides an
 603 attractive alternative to Laplace’s method when the model includes
 604 non-zero lapse rates, where the posterior may be less well approx-
 605 imated by a Gaussian. To compute mutual information from sam-
 606 ples, it is more convenient to use the expansion given in (eq. 15), so
 607 that it is expressed as the expected uncertainty reduction in entropy
 608 of the response \mathbf{y} , instead of a reduction in the posterior entropy.
 609 This will make it straightforward to approximate integrals needed
 610 for mutual information by Monte Carlo integrals involving sums
 611 over samples. Also note that we are back in the full parameter
 612 space; we no longer treat the lapse parameters separately, as we
 613 did for the Laplace-based infomax.

614 Given a set of posterior samples $\{\boldsymbol{\theta}_m\}$ from $p(\boldsymbol{\theta}|\mathcal{D}_t)$, the poste-
 615 rior distribution at time t , we can evaluate the mutual information
 616 using sums over “potential” terms that we denote by

$$617 L_{jm}(\mathbf{x}) \equiv p(y_j = 1|\mathbf{x}, \boldsymbol{\theta}_m). \quad (21)$$

618 This allows us to evaluate the conditional response entropy as

$$619 H_t(\mathbf{y}|\boldsymbol{\theta}; \mathbf{x}) \approx -\frac{1}{M} \sum_{j,m} L_{jm}(\mathbf{x}) \log L_{jm}(\mathbf{x}), \quad (22)$$

620 and the marginal response entropy as

$$621 H_t(\mathbf{y}; \mathbf{x}) \approx -\sum_j \left(\frac{1}{M} \sum_m L_{jm}(\mathbf{x}) \right) \log \left(\frac{1}{M} \sum_m L_{jm}(\mathbf{x}) \right), \quad (23)$$

622 where we have evaluated the posterior-predictive distribution as

$$623 p(y_j = 1|\mathbf{x}, \mathcal{D}_t) \approx \frac{1}{M} \sum_m L_{jm}(\mathbf{x}). \quad (24)$$

624 Putting together these terms, the mutual information can be evalu-
 625 ated as

$$626 I_t(\boldsymbol{\theta}; \mathbf{y}|\mathbf{x}) = -\frac{1}{M} \sum_{j,m} L_{jm}(\mathbf{x}) \log \frac{L_{jm}(\mathbf{x})}{\sum_{m'} L_{jm'}(\mathbf{x})/M}, \quad (25)$$

627 which is straightforward to evaluate for a set of candidate stimuli
628 $\{\mathbf{x}\}$. The computational cost of this approach is therefore linear
629 in the number of samples, and the primary concern is the cost of
630 obtaining a representative sample from the posterior.

631 Results

632 We consider two approaches for testing the performance of our pro-
633 posed stimulus-selection algorithms, one using simulated data, and
634 a second using an offline analysis of data from real psychophysical
635 experiments.

636 **Simulated experiments.** We first tested the performance of
637 our algorithms using simulated data from a fixed psychophysical
638 observer model. In these simulations, a stimulus \mathbf{x} was selected on
639 each trial and the observer’s response \mathbf{y} was sampled from a “true”
640 psychometric function, $p_{\text{true}}(\mathbf{y}|\mathbf{x}) = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_{\text{true}})$.

641 We considered psychophysical models defined on a continuous
642 2-dimensional stimulus space with 4 discrete response alternatives
643 for every trial, corresponding to the problem of estimating the di-
644 rection of 2D stimulus moving along one of the four cardinal di-
645 rections (up, down, left, right). We computed expected informa-
646 tion gain over a set of discrete stimulus values corresponding to
647 21×21 square grid (Fig. 5A). The stimulus plane is colored in
648 Fig. 5A, to indicate the most likely response (one of the four alter-
649 natives) in each stimulus region. Lapse probabilities λ_{c_i} were set
650 to either zero (the “lapse-free” case), or a constant value of 0.05,
651 resulting in a total lapse probability of $\lambda = 0.2$ across the four
652 choices (Fig. 5B). We compared performance of our adaptive algo-
653 rithms with a method that selected a stimulus uniformly at random
654 from the grid on each trial. We observed that the adaptive methods
655 tended to sample more stimuli near the boundaries between colored
656 regions on the stimulus space (Fig. 5C), which led to more efficient
657 estimates of the PF compared to the uniform stimulus selection ap-
658 proach (Fig. 5D). We also confirmed that the posterior entropy of
659 the inferred parameters decrease more rapidly with our adaptive
660 stimulus sampling algorithms, in all cases (Fig. 5E-F). This was

661 expected because our algorithms explicitly attempt to minimize the
662 posterior entropy, by maximizing the mutual information.

663 For each true model, we compared the performances of four dif-
664 ferent adaptive methods (Fig. 6A-B), defined by performing infer-
665 ence with MAP or MCMC, and assuming lapse rate to be fixed
666 at zero or including a non-zero lapse parameters. Each of these
667 inference methods was also applied to data selected according to
668 a uniform stimulus selection algorithm. We quantified perfor-
669 mance using the mean-squared error (MSE) between the true re-
670 sponse probabilities $p_{ij} = p(y = j|\mathbf{x}_i, \boldsymbol{\theta}_{\text{true}})$ and the estimated
671 probabilities \hat{p}_{ij} over the 21×21 grid of stimulus locations $\{\mathbf{x}_i\}$
672 and the 4 possible responses $\{j\}$. For MAP-based inference, es-
673 timated probabilities were given by $\hat{p}_{ij} = p(y = j|\mathbf{x}_i, \hat{\boldsymbol{\theta}}_{\text{MAP}})$.
674 For the MCMC-based inference, probabilities were given by the
675 predictive distribution, evaluated using an average over samples:
676 $\hat{p}_{ij} = \frac{1}{M} \sum_m p(y = j|\mathbf{x}_i, \boldsymbol{\theta}_m)$, where $\{\boldsymbol{\theta}_m\}$ represent samples
677 from the posterior.

678 When the true model was lapse-free (Fig. 6A), lapse-free and
679 lapse-aware inference methods performed similarly, indicating that
680 there was minimal cost to incorporating parameters governing
681 lapse when lapses were absent. Under all inference methods, in-
682 fomax stimulus selection outperformed uniform stimulus selec-
683 tion by a substantial margin. For example, infomax algorithms
684 achieved in 50 – 60 trials the error levels that their uniform-
685 stimulus-selection counterparts required 100 trials to achieve.

686 By contrast, when the true model had a non-zero lapse rate
687 (Fig. 6B), adaptive stimulus selection algorithms based on the
688 lapse-free model failed to select optimal stimuli, performing even
689 worse than uniform stimulus selection algorithms. This empha-
690 sizes the impact of model mismatch in adaptive methods, and the
691 importance of a realistic psychometric model. When lapse-aware
692 models were used for inference, on the other hand, both Laplace-
693 based and MCMC-based adaptive stimulus selection algorithms
694 achieved a significant speedup compared to uniform stimulus se-
695 lection, while MCMC-based adaptive algorithm performed bet-
696 ter. This shows that the MCMC-based infomax stimulus selec-

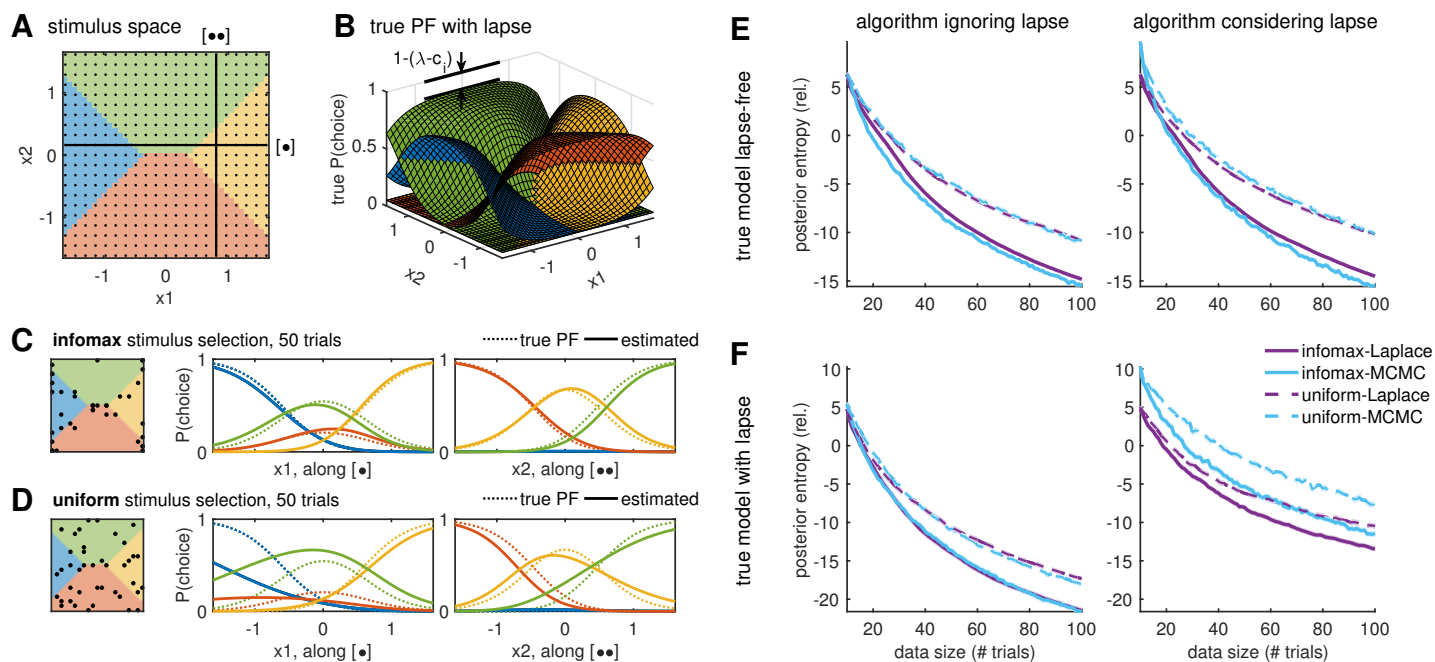


Figure 5: The simulated experiment. (A) At each trial, a stimulus was selected from a 2D stimulus plane with a 21×21 grid. The two lines, running along x_1 and x_2 respectively, indicate the cross-sections used in C and D below. Colors indicate the most likely response in the respective stimulus regime, according to the true PF shown in B, with a consistent color code. (B) Given each stimulus, a simulated response was drawn from a true model with 4 alternatives. Shown here is the model with lapse, characterized by a non-deterministic choice (i.e., the choice probability does not approach 0 or 1) even at an easy stimulus, far from the choice boundaries. (C-D) Examples of Laplace-approximation-based inference results after 50 trials, where stimuli was selected either using our adaptive infomax method (C) or uniformly (D), as shown on left. In both cases, the true model was lapse-free, and the algorithm assumed that lapse was fixed at zero. The two sets of curves show the cross-sections of the true PF (dotted lines) and the estimated PF (solid lines), along the two lines marked in A, after sampling these stimuli. (E-F) Traces of posterior entropy from simulated experiments, averaged over 100 runs each. The true model for simulation was either (E) lapse-free, or (F) with a finite lapse rate of $\lambda = 0.2$, with a uniform lapse scenario $c_i = 1/4$ for each outcome $i = 1, 2, 3, 4$. In algorithms considering lapse (panels on the right), the shift in posterior entropy is due to the use of partial covariance (with respect to weight) in the case of Laplace approximation. The algorithm either used the classical MNL model that assumes zero lapse (left column), or our extended model that considers lapse (right column). Average performances of adaptive and uniform stimulus selection algorithms are plotted in solid and dashed lines, respectively; Laplace-based and MCMC-based algorithms are plotted in purple and cyan. The lighter lines show standard error intervals over 100 runs, which are very narrow. All sampling-based algorithms used the semi-adaptive MCMC with chain length $M = 1000$.

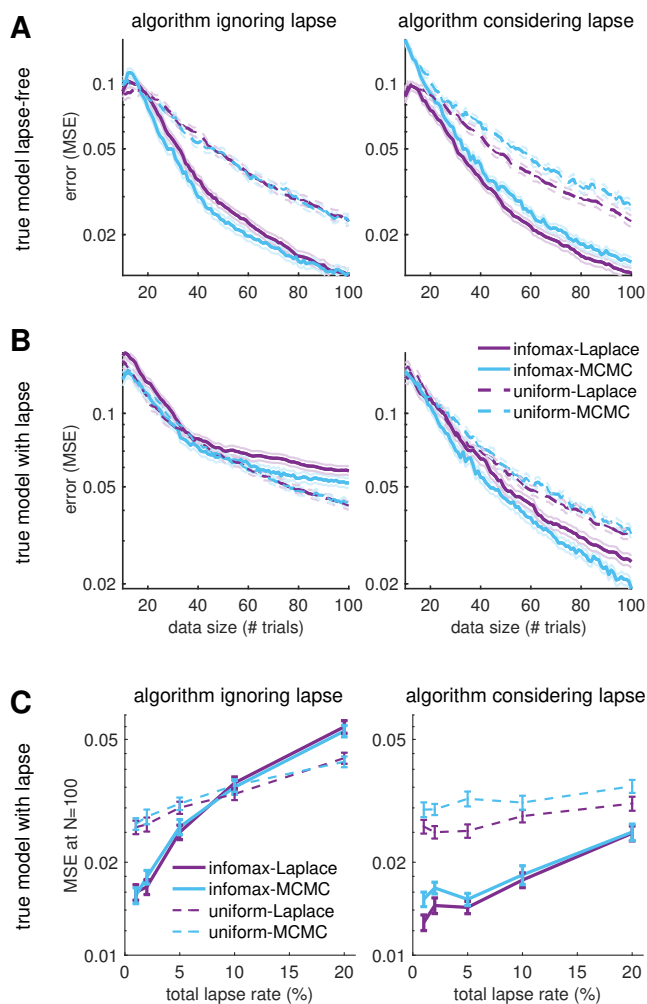


Figure 6: The simulated experiment, continued. We show results from the same set of simulated experiments as in Fig. 5. (A-B) Traces of the mean-square error (MSE), where the true model was either (A) lapse-free, or (B) with a total lapse rate of $\lambda = 0.2$, uniformly distributed to each outcome. Standard error intervals are plotted in lighter lines as in Fig. 5E-F. (C) Effect of lapse, tested by adding varying total lapse rates λ . Shown are the MSE after $N = 100$ trials of each stimulus selection algorithm, equivalent to the endpoints in B. Error bars indicate the standard error over 100 runs, equivalent to the lighter-line intervals in the above panels.

697 tion method can provide an efficient and robust platform for adap-
 698 tive experiments with realistic models. When the true behavior
 699 had lapses, the MCMC-based adaptive stimulus selection algo-
 700 rithm with the lapse-aware model automatically included “easy”
 701 trials, which provide maximal information about lapse probabili-
 702 ties. These easy trials are typically in the periphery of the stimulus

space (strong-stimulus regimes, referred to as “asymptotic perfor-
 mance intensity” in Prins (2012)).

However, that the effect of model mismatch due to non-zero lapse only becomes problematic at high enough lapse rate; in the simulation shown in Fig. 5F and Fig. 6B, we used a high lapse rate of $\lambda = 0.2$ which is more typical in the case of less sophisticated animals such as rodents (see for example Scott, Constantino-
 ple, Erlich, Tank, and Brody (2015)). With lapse rates more typical in well-designed human psychophysics tasks ($\lambda \lesssim 0.05$; see for example Wichmann and Hill (2001a, 2001b)), infomax algorithms still tend to perform better than uniform sampling algorithms (Fig. 6C).

Finally, we measured the computation time per trial required by our adaptive stimulus selection algorithms on a personal desktop with an Intel i7 processor. With the Laplace-based algorithm, the major computational bottleneck is the parameter space integration in the infomax calculation, which scales directly with the model complexity. We could easily achieve tens-of-milliseconds trials in the case of the simple 2AFC task, and sub-second trials with 2-dimensional stimuli and 4-alternative responses, as used in the current set of simulations (Fig. 7A-B). With the MCMC-based algorithm, the time-per-trial in the sampling-based method is limited by the number of samples in each MCMC chain, M , rather than by the model complexity. Using the standard implementation for the Metropolis-Hastings sampler in Matlab, a time-per-trial of ~ 0.1 seconds was achieved with chains shorter than $M \lesssim 200$ (Fig. 7C-D, top panels). This length of $M \approx 200$ was good enough to represent the posterior distributions for our simulated examples (Fig. 7C-D, bottom panels), although we note that longer chains are required to sample a more complex posterior distribution, and this particular length M should not be taken as the benchmark in general.

Optimal re-ordering of real dataset. A second approach for testing the performance of our methods is to perform an off-line analysis of data from real psychophysical experiments. Here we take an existing dataset and use our methods to re-order the trials so

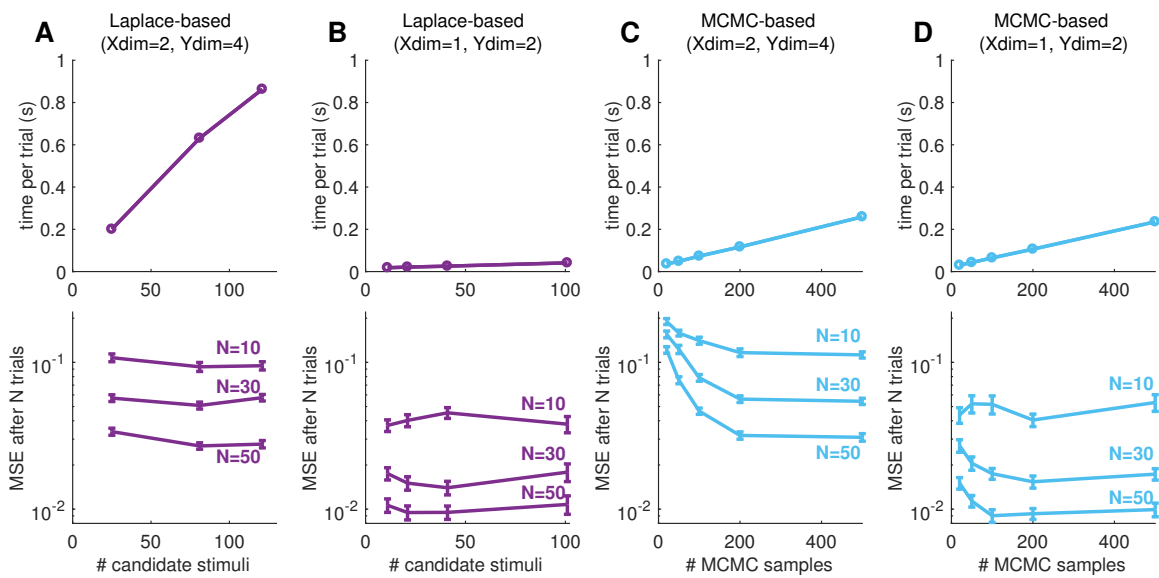


Figure 7: Computation time and accuracy. (A-B) The computation times for the Laplace-based algorithms grow linearly with the number of candidate stimulus points, as shown on the top panels, because one needs to perform a numerical integration to compute the expected utility of each stimulus. In general, there is a tradeoff between cost (computation time) and accuracy (inversely related to the estimation error). The bottom panels show the mean-square error of the estimated PF, calculated after completing a sequence of N trials, where the 10 initial trials were selected at regular intervals, and the following trials were selected under our adaptive algorithm. Error estimates were averaged over 100 independent sequences. Error bars indicate the standard errors. The true model used were the same as either (A) in Fig. 5, with 2-dimensional stimulus and 4-alternative response, described by 9 parameters; or (B) in Fig. 3, with 1-dimensional stimulus and binary response, with only 2 parameters (slope and threshold). Different rate at which the computation time increases under the two model reflects the different complexity of numerical quadrature involved. We used lapse-free algorithms in all cases in this example. (C-D) We similarly tested the MCMC-based algorithms using the two models as in panels A-B. In this case, the computation times (top panels) grow linearly the number of samples in each MCMC chain, and are not sensitive to the dimensionality of the parameter space. On the other hand, the estimation error plots (bottom panels) suggest that a high-dimensional model requires more samples for accurate inference.

739 that the most-informative stimuli are selected first (also see Lewi,
 740 Schneider, Woolley, and Paninski (2011) for a similar approach).
 741 To obtain a re-ordering, we iteratively apply our algorithm to the
 742 stimuli shown during the experiment. On each trial, we use our
 743 adaptive algorithm to select the optimal stimulus from the set of
 744 stimuli $\{x_i\}$ not yet incorporated into the model. This selection
 745 takes place without access to the actual responses $\{y_i\}$. We then
 746 update the posterior using the stimulus x_i and the response y_i it ac-
 747 tually elicited during the experiment, then proceed to the next trial.
 748 We can then ask whether adding the data according to the proposed
 749 re-ordering would have led to faster narrowing of the posterior dis-
 750 tribution than other orderings.

751 To perform this analysis, we used a dataset from macaque

752 monkeys performing a four-alternative motion discrimination task
 753 (Churchland, Kiani, & Shadlen, 2008). Monkeys were trained to
 754 observe a motion stimulus with dots moving in one of the four car-
 755 dinal directions, and report this direction of motion with an eye
 756 movement. The difficulty of the task was controlled by varying the
 757 fraction of coherently moving dots on each trial, with the remain-
 758 ing dots appearing randomly (Fig. 8A). Each moving-dot stimulus
 759 in this experiment could be represented as a two-dimensional vec-
 760 tor, where the direction of the vector is the direction of the mean
 761 movement of the dots, and the amplitude of the vector is given by
 762 the fraction of coherently moving dots (a number between 0 and
 763 1). Each stimulus presented in the the experiment was aligned with
 764 either one of the two cardinal axes of the stimulus plane (Fig. 8B).

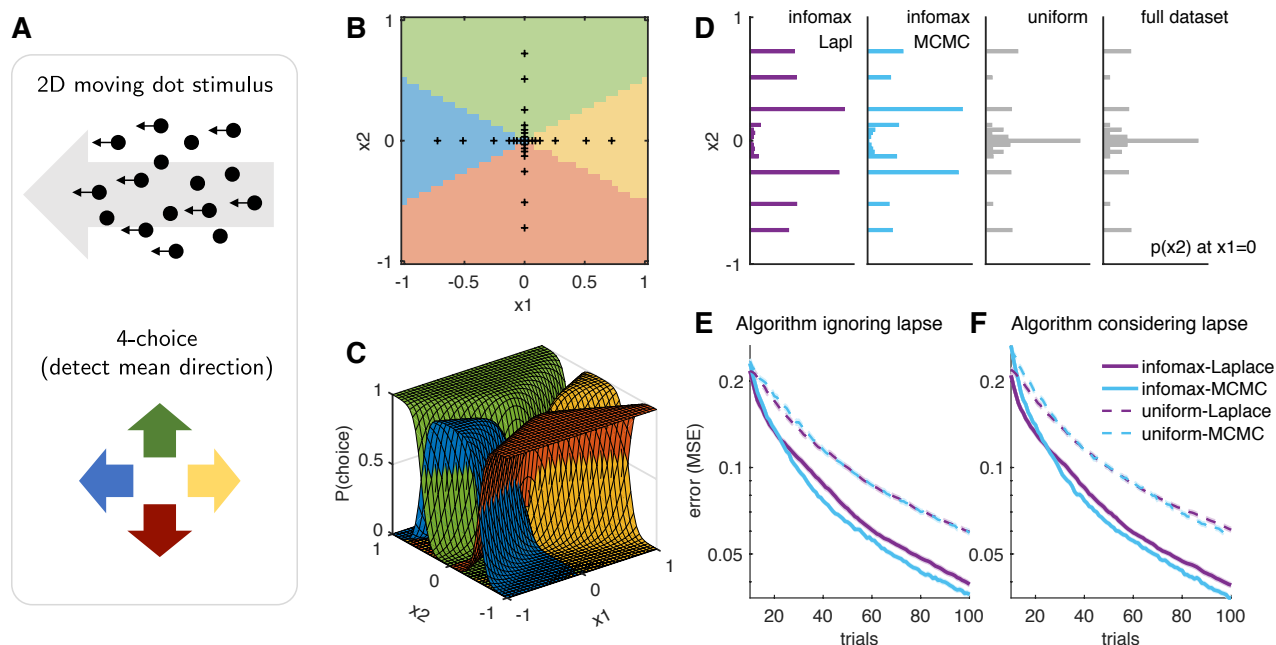


Figure 8: Optimal re-ordering of a real monkey dataset. (A) The psychometric task consisted of a 2D stimulus presented as moving dots, characterized by a coherence and a mean direction of movement, and a 4-alternative response. The four choices are color coded consistently in A-C in this figure. (B) The axes-only stimulus space of the original dataset, with 15 fixed stimuli along each axis. Colors indicate the most likely response in the respective stimulus regime according to the best estimate of the PF. (C) The best estimate of the PF of monkeys in this task, inferred from all observations in the dataset. (D) Stimuli selection in the first $N = 100$ trials during the re-ordering experiment, under the inference method that ignores lapse. Shown are histograms of x_2 along one of the axes, $x_1 = 0$, averaged over 100 independent runs in each case. (E-F) Error traces under different algorithms, averaged over 100 runs. Both Laplace-based (purple) and MCMC-based (cyan; with $M = 1000$) algorithms achieve significant speedups over uniform sampling. Because the monkeys were almost lapse-free in this task, inference methods that ignore lapse (E) and consider lapse (F) performed similarly. Standard error intervals over 100 runs are shown in lighter lines, but are very narrow.

765 The PF for this dataset consists of a set of four 2D curves, where
 766 each curve specifies the probability of choosing a particular direc-
 767 tion as a function of location in the 2D stimulus plane (Fig. 8C).

768 This monkey dataset contained more than 10,000 total obser-
 769 vations at 29 distinct stimulus conditions, accumulating more than
 770 300 observations per stimulus. This multiplicity of observations
 771 per stimulus ensured that the posterior distribution given the full
 772 dataset was narrow enough that it could be considered to provide a
 773 “ground truth” psychometric function against which the inferences
 774 based on the re-ordering experiment could be compared.

775 The first 100 stimuli selected by the infomax algorithms had
 776 noticeably different statistics than the full dataset or its uniform
 777 sub-sampling (the first $N = 100$ trials under uniform sampling).

778 On the other hand, the sets of stimuli selected by both MAP-
 779 based and MCMC-based infomax algorithms were similar. Fig. 8D
 780 shows the histogram of stimulus component along one of the axes,
 781 $p(x_2 | x_1 = 0)$, from the first $N = 100$ trials, averaged over 100
 782 independent runs under each stimulus selection algorithm using the
 783 lapse-free model.

784 Because the true PF was unknown, we compared the perfor-
 785 mance of each algorithm to an estimate of the PF from the entire
 786 dataset. When using the MAP algorithm, the full-dataset PF was
 787 given by $p_{ij} = p(y = j | \mathbf{x}_i, \hat{\theta}_{\text{full}})$, evaluated at the MAP estimate
 788 of the log posterior, $\hat{\theta}_{\text{full}} = \arg\max_{\theta} \log p(\theta | \mathcal{D}_{\text{full}})$, given the full
 789 dataset $\mathcal{D}_{\text{full}}$. For the MCMC algorithm, the full-dataset PF was
 790 computed by $p_{ij} \approx \frac{1}{M} \sum_m p(y = j | \mathbf{x}_i, \theta_m)$, where the MCMC

chain $\{\theta_m\} \sim \log p(\theta | \mathcal{D}_{\text{full}})$ sampled the log posterior given the full dataset. The re-ordering test on the monkey dataset showed that our adaptive stimulus sampling algorithms were able to infer the PF to a given accuracy in a smaller number of observations, compared to a uniform sampling algorithm (Fig. 8E-F). In other words, data collection could have been faster with an optimal re-ordering of the experimental procedure.

Exploiting the full stimulus space. In the experimental dataset considered in the previous section, the motion stimuli were restricted to points along the cardinal axes of the 2D stimulus plane (Fig. 8B) (Churchland et al., 2008). In some experimental settings, however, the psychometric functions of interest may lack identifiable axes of alignment or may exhibit asymmetries in shape or orientation. Here we show that in such cases, adaptive stimulus selection methods can benefit from the ability to select points from the full space of possible stimuli.

We performed experiments with a simulated observer governed by the lapse-free psychometric function estimated from the macaque monkey dataset (Fig. 8C). This psychometric function was either aligned to the original stimulus axes (Fig. 9A-B) or rotated counter-clockwise by 45 degrees (Fig. 9C). We tested the performance of adaptive stimulus selection using the Laplace infomax algorithm, with stimuli restricted to points along the cardinal axes (Fig. 9A), or allowed to a grid of points in the full 2D stimulus plane (Fig. 9B-C).

The simulated experiment indeed closely resembled the results of our dataset re-ordering test in terms of the statistics of adaptively selected stimuli (compare Fig. 9A to the purple histogram in Fig. 8D). With the full 2D stimulus space aligned to the cardinal axes, on the other hand, our adaptive infomax algorithm detected and sampled more stimuli near the boundaries between colored regions in the stimulus plane, which were usually not on the cardinal axes (Fig. 9B). Finally, we also observed that this automatic exploitation of the stimulus space was not limited by the lack of alignment between the PF and the stimulus axes; our adaptive infomax algorithm was just as effective in detecting and sampling the

boundaries between stimulus regions in the case of the unaligned PF (Fig. 9C).

The error traces in Fig. 9D show that we can infer the PF at a given accuracy in an even fewer number of observations using our adaptive algorithm on the full 2D stimulus plane (orange curves), compared to the cardinal-axes design (black curves). It also confirms that we can infer the PF accurately and effectively with an unaligned stimulus space (red curves), as well as with an aligned stimulus space. For comparison purposes, all errors were calculated over the same 2D stimulus grid, even when the stimulus selection was from the cardinal axes. (This had negligible effects on the resulting error values: compare the black curves in Fig. 9D and the purple curves in Fig. 8E.)

Discussion

We developed effective Bayesian adaptive stimulus selection algorithms for inferring psychometric functions, with an objective of maximizing the expected informativeness of each stimulus. The algorithms select an optimal stimulus adaptively in each trial, based on the posterior distribution of model parameters inferred from the accumulating set of past observations.

We emphasized that in psychometric experiments, especially with animals, it is crucial to use models that can account for the non-ideal yet common behaviors, such as omission (no response; an additional possibility for the outcome) or lapse (resulting in a random, stimulus-independent response). Specifically, we constructed a hierarchical extension of a multinomial logistic (MNL) model that incorporates both omission and lapse. Although we did not apply these additional features to real data, we performed simulated experiments to investigate their impacts on the accurate inference of psychometric functions. To ensure applicability of the extended model in real-time closed-loop adaptive stimulus selection algorithms, we also developed efficient methods for inferring the posterior distribution of the model parameters, with approximations specifically suited for sequential experiments.

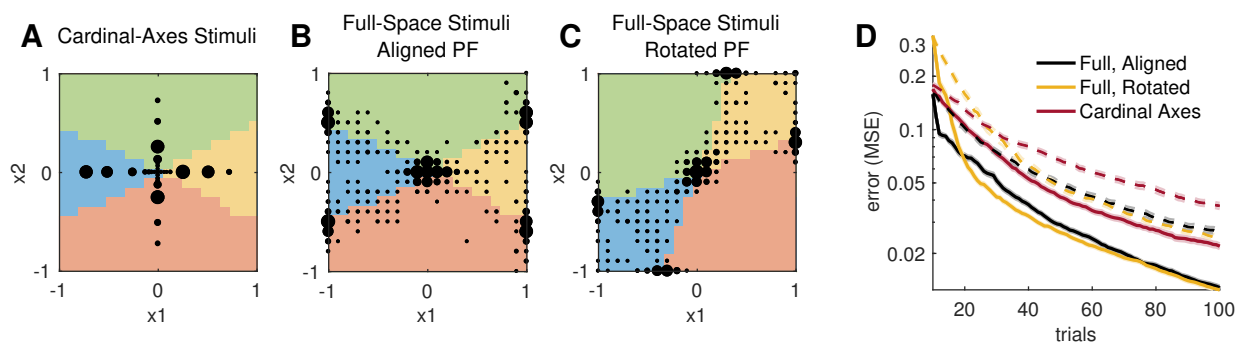


Figure 9: Design of multi-dimensional stimulus space. (A-C) Three different stimulus space designs were used in a simulated psychometric experiment. Responses were simulated according to fixed lapse-free PFs, matched to our best estimate of the monkey PF (Fig. 8C). Stimuli were selected within the respective stimulus spaces, (A) the cardinal-axes design, as in the original experiment; (B) full stimulus plane, with the PF aligned to the cardinal axes of the original stimulus space; (C) full stimulus plane, with rotated PF. The black dots in A-C indicate which stimuli were sampled by the Laplace-based infomax algorithm during the first $N = 100$ trials of simulation, where the dot size is proportional to the number of trials in which each stimulus was selected (averaged over 20 independent runs, and excluding the 10 fixed initial stimuli). (D) The corresponding error traces, under infomax (solid lines) or uniform (dashed lines) stimulus selection, averaged over 100 runs respectively. Colors indicate the three stimulus space designs, as shown in A-C. Standard error intervals over 100 runs are shown in lighter lines.

861 **Advantages of adaptive stimulus selection.** We observed
 862 two important advantages of using Bayesian adaptive stimulus se-
 863 lection methods in psychometric experiments. First, we showed
 864 that our adaptive stimulus selection algorithms achieved signifi-
 865 cant speed-ups in learning time (number of measurements), both
 866 on simulated data and in re-ordering test of a real experimental
 867 dataset, with and without lapse in the underlying behavior. Import-
 868 antly, the success of the algorithm depends heavily on the use of
 869 the correct model family; for example, adaptive stimulus selection
 870 fails when a classical (lapse-ignorant) model was used to measure
 871 behavior with a finite lapse rate. Based on the simulation results,
 872 it seems good practice to always use the lapse-aware model unless
 873 the behavior under study is known to be completely lapse-free, al-
 874 though it should be checked that the addition of the lapse param-
 875 eters does not make the inference problem intractable, given the
 876 constraints of the specific experiments. (One way to check this is
 877 using a simulated experiment, where lapse is added to the psycho-
 878 metric function inferred by lapse-free model; similarly to what we
 879 did in this paper.) The computational cost for incorporating lapses
 880 amounts to having k additional parameters to sample, one per each
 881 available choice, which is independent from the dimensionality of

the stimulus space.

882
 883 Second, we demonstrated that our adaptive stimulus selection
 884 study has implications on the optimization of the experimental de-
 885 signs more generally. Contrary to the conventional practice of ac-
 886 cumulating repeated observations at a small set of fixed stimuli, we
 887 suggest that the (potentially high-dimensional) stimulus space can
 888 be exploited more efficiently using our Bayesian adaptive stimulus
 889 selection algorithm. Specifically, the adaptive stimulus selection
 890 algorithm can automatically detect the structure of the stimulus
 891 space (with respect to the psychometric function) as part of the
 892 process. We also showed that there are benefits of using the full
 893 stimulus space even when the PF is aligned to the cardinal axes of
 894 the stimulus space.

895 **Comparison of the two algorithms.** Our adaptive stimulus
 896 selection algorithms were developed based on two methods for ef-
 897 fective posterior inference: one based on local Gaussian approxi-
 898 mation (Laplace approximation) of the posterior, and another based
 899 on MCMC sampling. The well-studied analytical method based
 900 on the Laplace approximation is fast and effective in simple cases,
 901 but becomes heavier in the case of more complicated PFs, because

902 the computational bottleneck is the numerical integration over the
903 parameter space that needs to be performed separately for each
904 candidate stimulus. In the case of sampling-based methods, on
905 the other hand, the computational speed is constrained by the num-
906 ber of MCMC samples used to approximate the posterior distribu-
907 tion, but not directly by the number of parameters or the number
908 of candidate stimuli. In general, however, accurately inferring a
909 higher-dimensional posterior distribution requires more samples,
910 and therefore a longer computation time. We note that our semi-
911 adaptive turning algorithm helps with the cost-accuracy tradeoff
912 by optimizing the sampling accuracy in a given number of sam-
913 ples, without human intervention. although it does not reduce the
914 computation time itself.

915 To summarize, when the PF under study is low-dimensional and
916 well-described by the multinomial logistic model, for example in
917 a 2AFC study with human subjects, Laplace-based approach pro-
918 vides a lightweight and elegant approach. But if the PF is higher-
919 dimensional or deviates significantly from the ideal model (e.g.,
920 large lapse), MCMC sampling provides a flexible and affordable
921 solution. Results suggest that our MCMC-based algorithm will be
922 applicable to most animal psychometric experiments, as the model
923 complexities are not expected to significantly exceed our simulated
924 example. However, one should always make sure that the number
925 of MCMC samples being used is sufficient to sample the posterior
926 distribution under study.

927 **Limitations and Open Problems.** One potential drawback
928 of adaptive experiments is the undesired possibility that the psy-
929 chometric function of the observer might adapt to the distribution
930 of stimuli presented during the experiments. If this is the case,
931 the system under measurement would no longer be stationary, nor
932 independent of the experimental design, profoundly altering the
933 problem one should try to solve. The usual assumption in psycho-
934 metric experiments is that well trained observers exhibit stationary
935 behavior on the timescale of an experiment; under this assumption,
936 the order of data collection cannot bias inference [MacKay \(1992\)](#).
937 However, the empirical validity of this claim remains a topic for

future research. 938

939 One approach for mitigating non-stationarity is to add regressors
940 to account for the history dependence of psychophysical behavior.
941 Recent work has shown that extending a psychophysical model to
942 incorporate past rewards ([Bak et al., 2016](#); [Busse et al., 2011](#); [Corrado, Sugrue, Seung, & Newsome, 2005](#); [Lau & Glimcher, 2005](#)),
943 past stimuli ([Akrami, Kopec, Diamond, & Brody, 2018](#)) or the full
944 stimulus-response history ([Fründ, Wichmann, & Macke, 2014](#)) can
945 provide a more accurate description of the factors influencing re-
946 sponses on a trial-by-trial basis. 947

948 Our work leaves open a variety of directions for future research.
949 One simple idea is to re-analyze old datasets under the multinomial
950 response model with omissions included as a separate response cat-
951 egory; this will reveal whether omissions exhibit stimulus depen-
952 dence (e.g., occurring more often on difficult trials), and provide
953 greater insight into the factors influencing psychophysical behavior
954 on single trials. Another set of directions is to extend the multino-
955 mial logistic observer model to obtain a more accurate or more
956 flexible model of psychophysical behavior; particular directions
957 include models with nonlinear stimulus dependencies or interac-
958 tion terms ([Cowley, Williamson, Clemens, Smith, & Byron, 2017](#);
959 [DiMattina & Zhang, 2011](#); [Hyafil & Moreno-Bote, 2017](#); [Neri & Heeger, 2002](#)),
960 models with output nonlinearities other than the
961 logistic ([Kontsevich & Tyler, 1999](#); [Schütt et al., 2016](#); [A. B. Watson, 2017](#);
962 [A. B. Watson & Pelli, 1983](#)), or models that capture
963 overdispersion, e.g., due to non-stationarities of the observer, via a
964 hierarchical prior ([Schütt et al., 2016](#)). In general, such extensions
965 will be much easier to implement with the MCMC-based inference
966 method, due to the fact that it does not rely on gradients or Hessians
967 of a particular parametrization of log-likelihood. Finally, it may be
968 useful to consider the same observer model under optimality cri-
969 teria other than mutual information — recent work has shown that
970 infomax methods do not necessarily attain optimal performance
971 according to alternate metrics (e.g., mean squared error, [I. M. Park and Pillow \(2017\)](#);
972 [M. Park et al. \(2014\)](#)) — or using non-greedy
973 selection criteria that optimize stimulus selection based on a time

974 horizon longer than the next trial (Kim et al., 2017; King-Smith et
975 al., 1994).

976 Code availability

977 A Matlab implementation of our methods is available online at
978 <https://github.com/pillowlab/adaptivePsychophysicsToolbox>.

979 Acknowledgements

980 We thank Anne Churchland for sharing the monkey data. JHB was
981 supported by the Samsung Scholarship for the study at Princeton.
982 JWP was supported by grants from the McKnight Foundation, Si-
983 mons Collaboration on the Global Brain (SCGB AWD1004351)
984 and the NSF CAREER Award (IIS-1150186). Computational work
985 was performed using resources at Princeton University and the
986 KIAS Center for Advanced Computing.

987 References

988 Akrami, A., Kopec, C. D., Diamond, M. E., & Brody, C. D. (2018).
989 Posterior parietal cortex represents sensory history and me-
990 diates its effects on behaviour. *Nature*, *554*(7692), 368–372.
991 doi: 10.1038/nature25510

992 Bak, J. H., Choi, J. Y., Akrami, A., Witten, I. B., & Pillow, J. W.
993 (2016). Adaptive optimal training of animal behavior. In
994 *Advances in neural information processing systems* 29 (pp.
995 1947–1955).

996 Barthelmé, S., & Mamassian, P. (2008). A flexi-
997 ble bayesian method for adaptive measurement in psy-
998 chophysics. *arXiv:0809.0387*, 1–28.

999 Bishop, C. M. (2006). *Pattern recognition and machine learning*.
1000 Springer New York.

1001 Busse, L., Ayaz, A., Dhruv, N. T., Katzner, S., Saleem, A. B.,
1002 Scholvinck, M. L., ... Carandini, M. (2011). The de-
1003 tection of visual contrast in the behaving mouse. *Jour-*

nal of Neuroscience, *31*(31), 11351–11361. doi: 10.1523/
JNEUROSCI.6689-10.2011

Carandini, M., & Churchland, A. K. (2013). Probing perceptual
decisions in rodents. *Nature Neuroscience*, *16*(7), 824–831.
doi: 10.1038/nn.3410

Cavagnaro, D. R., Myung, J. I., Pitt, M. A., & Kujala, J. V. (2010).
Adaptive design optimization: a mutual information-based
approach to model discrimination in cognitive science. *Neu-
ral computation*, *22*(4), 887–905. doi: 10.1162/neco.2009
.02-09-959

Chaloner, K., & Larntz, K. (1989). Optimal logistic Bayesian
design applied to logistic regression experiments. *Journal
of Statistical Planning and Inference*, *21*, 191–208. doi: 10
.1016/0378-3758(89)90004-9

Chaloner, K., & Verdinelli, I. (1995). Bayesian experimental de-
sign: a review. *Statistical Science*, *10*, 273–304.

Churchland, A. K., Kiani, R., & Shadlen, M. N. (2008). Decision-
making with multiple alternatives. *Nature Neuroscience*,
11(6), 693–702. doi: 10.1038/nn.2123

Corrado, G. S., Sugrue, L. P., Seung, H. S., & Newsome, W. T.
(2005). Linear-nonlinear-poisson models of primate choice
dynamics. *Journal of the Experimental Analysis of Behavior*,
84(3), 581–617. doi: 10.1901/jeab.2005.23-05

Cowley, B., Williamson, R., Clemens, K., Smith, M., & Byron,
M. Y. (2017). Adaptive stimulus selection for optimizing
neural population responses. In *Advances in neural infor-
mation processing systems* (pp. 1395–1405).

DiMattina, C. (2015). Fast adaptive estimation of multidimen-
sional psychometric functions. *Journal of Vision*, *15*(9), 5.
doi: 10.1167/15.9.5

DiMattina, C., & Zhang, K. (2011). Active data collection
for efficient estimation and comparison of nonlinear neu-
ral models. *Neural Computation*, *23*(9), 2242–2288. doi:
10.1162/NECO_a_00167

- 1038 Fründ, I., Wichmann, F. A., & Macke, J. H. (2014). Quantifying
1039 the effect of intertrial dependence on perceptual decisions.
1040 *Journal of vision*, *14*(7), 1–16. doi: 10.1167/14.7.9
- 1041 Gardner, J. R., Song, X., Weinberger, K. Q., Barbour, D., & Cun-
1042 ningham, J. P. (2015). Psychophysical detection testing with
1043 bayesian active learning. In *Proceedings of the thirty-first*
1044 *conference on uncertainty in artificial intelligence* (pp. 286–
1045 297). AUAI Press.
- 1046 Gelman, A., Roberts, G., & Gilks, W. (1996). Efficient metropolis
1047 jumping rules. *Bayesian statistics*, *5*, 599–607.
- 1048 Glonek, G., & McCullagh, P. (1995). Multivariate Logistic
1049 Models. *Journal of the Royal Statistical Society, Series B*
1050 *(Methodological)*, *57*(3), 533–546.
- 1051 Haario, H., Saksman, E., & Tamminen, J. (2001). An adaptive
1052 metropolis algorithm. *Bernoulli*, *7*(2), 223–242. doi: 10
1053 .2307/3318737
- 1054 Heiss, F., & Winschel, V. (2008). Likelihood approximation by
1055 numerical integration on sparse grids. *Journal of Economet-*
1056 *rics*, *144*(1), 62–80. doi: 10.1016/j.jeconom.2007.12.004
- 1057 Henderson, H. V., & Searle, S. R. (1981). On deriving the inverse
1058 of a sum of matrices. *SIAM Review*, *23*(1), 53–60. doi:
1059 10.1137/1023004
- 1060 Higham, N. J. (1988). Computing a nearest symmetric positive
1061 semidefinite matrix. *Linear Algebra and Its Applications*,
1062 *103*(C), 103–118. doi: 10.1016/0024-3795(88)90223-6
- 1063 Hyafil, A., & Moreno-Bote, R. (2017). Breaking down hierarchies
1064 of decision-making in primates. *eLife*, *6*.
- 1065 Kim, W., Pitt, M. A., Lu, Z., & Myung, J. I. (2017). Planning be-
1066 yond the next trial in adaptive experiments: A dynamic pro-
1067 gramming approach. *Cognitive Science*, *41*(8), 2234–2252.
1068 doi: 10.1111/cogs.12467
- 1069 Kim, W., Pitt, M. A., Lu, Z.-l., Steyvers, M., & Myung, J. I.
1070 (2014). A hierarchical adaptive approach to optimal ex-
1071 perimental design paradigm of adaptive design optimiza-
tion (ADO). *Neural computation*, *26*, 2465–2492. doi:
10.1162/NECO_a.00654
- King-Smith, P. E., Grigsby, S. S., Vingrys, A. J., Benes, S. C.,
& Supowit, A. (1994, apr). Efficient and unbiased modifi-
cations of the quest threshold method: Theory, simulations,
experimental evaluation and practical implementation. *Vi-
sion Research*, *34*(7), 885–912. doi: 10.1016/0042-6989(94)
90039-6
- Knoblauch, K., & Maloney, L. T. (2008). Estimating classification
images with generalized linear and additive models. *Journal
of Vision*, *8*(16), 10.1–1019. doi: 10.1167/8.16.10
- Kontsevich, L. L., & Tyler, C. W. (1999). Bayesian adaptive es-
timation of psychometric slope and threshold. *Vision Re-
search*, *39*(16), 2729–2737. doi: 10.1016/S0042-6989(98)
00285-5
- Kujala, J. V., & Lukka, T. J. (2006). Bayesian adaptive estimation:
The next dimension. *Journal of Mathematical Psychology*,
50(4), 369–389. doi: 10.1016/j.jmp.2005.12.005
- Kuss, M., Jäkel, F., & Wichmann, F. A. (2005). Bayesian inference
for psychometric functions. *Journal of Vision*, *5*(5), 478–
492. doi: 10.1167/5.5.8
- Lau, B., & Glimcher, P. W. (2005). Dynamic response-by-response
models of matching behavior in rhesus monkeys. *Journal of
the Experimental Analysis of Behavior*, *84*(3), 555–579. doi:
10.1901/jeab.2005.110-04
- Lesmes, L. A., Lu, Z.-L., Baek, J., & Albright, T. D. (2010).
Bayesian adaptive estimation of the contrast sensitivity func-
tion: the quick CSF method. *Journal of Vision*, *10*(3), 17.1–
21. doi: 10.1167/10.3.17
- Lesmes, L. A., Lu, Z.-L., Baek, J., Tran, N., Doshier, B., & Al-
bright, T. (2015). Developing bayesian adaptive meth-
ods for estimating sensitivity thresholds (d?) in yes-no and
forced-choice tasks. *Frontiers in Psychology*, *6*, 1070. doi:
10.3389/fpsyg.2015.01070
- Lewi, J., Butera, R., & Paninski, L. (2009). Sequential optimal de-

- 1107 sign of neurophysiology experiments. *Neural computation*, 1141
1108 *21*(3), 619–687. doi: 10.1162/neco.2008.08-07-594 1142
- 1109 Lewi, J., Schneider, D. M., Woolley, S. M. N., & Paninski, L. 1143
1110 (2011). Automating the design of informative sequences 1144
1111 of sensory stimuli. *Journal of Computational Neuroscience*, 1145
1112 *30*(1), 181–200. doi: 10.1007/s10827-010-0248-1 1146
- 1113 MacKay, D. J. C. (1992). Information-based objective functions 1147
1114 for active data selection. *Neural Computation*, *4*(4), 590– 1148
1115 604. doi: 10.1162/neco.1992.4.4.590 1149
- 1116 Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, 1150
1117 A. H., & Teller, E. (1953). Equation of state calculations by 1151
1118 fast computing machines. *The Journal of Chemical Physics*, 1152
1119 *21*(6), 1087. doi: 10.1063/1.1699114 1153
- 1120 Murray, R. F. (2011). Classification images: A review. *Journal of* 1154
1121 *Vision*, *11*(5). doi: 10.1167/11.5.2 1155
- 1122 Neri, P., & Heeger, D. J. (2002). Spatiotemporal mechanisms for 1156
1123 detecting and identifying image features in human vision. 1157
1124 *Nat Neurosci*, *5*(8), 812–816. doi: 10.1038/nn886 1158
- 1125 Paninski, L., Ahmadian, Y., Ferreira, D. G., Koyama, S., Rah- 1159
1126 nama Rad, K., Vidne, M., ... Wu, W. (2010). A new 1160
1127 look at state-space models for neural data. *Journal of Com-* 1161
1128 *putational Neuroscience*, *29*(1), 107–126. doi: 10.1007/ 1162
1129 s10827-009-0179-x 1163
- 1130 Park, I. M., & Pillow, J. W. (2017). Bayesian efficient coding. 1164
1131 *bioRxiv*, 178418. doi: 10.1101/178418 1165
- 1132 Park, M., Horwitz, G., & Pillow, J. W. (2011). Active learning 1166
1133 of neural response functions with Gaussian processes. In 1167
1134 *Advances in neural information processing systems 24* (pp. 1168
1135 2043–2051). 1169
- 1136 Park, M., & Pillow, J. W. (2012). Bayesian active learning with 1170
1137 localized priors for fast receptive field characterization. In 1171
1138 *Advances in neural information processing systems 25* (pp. 1172
1139 2357–2365). 1173
- 1140 Park, M., Weller, J. P., Horwitz, G. D., & Pillow, J. W. (2014). 1174
Bayesian active learning of neural firing rate maps with 1175
transformed gaussian process priors. *Neural Computation*,
26(8), 1519–1541.
- Pillow, J. W., Ahmadian, Y., & Paninski, L. (2011). Model-based
decoding, information estimation, and change-point detec-
tion techniques for multineuron spike trains. *Neural Com-*
putation, *23*(1), 1–45. doi: 10.1162/NECO_a.00058
- Pillow, J. W., & Park, M. (2016). Adaptive Bayesian methods for
closed-loop neurophysiology. In A. E. Hady (Ed.), *Closed*
loop neuroscience. Elsevier.
- Prins, N. (2012). The psychometric function: The lapse rate revis-
ited. *Journal of Vision*, *12*(6), 25–25. doi: 10.1167/12.6.25
- Prins, N. (2013). The psi-marginal adaptive method: How to give
nuisance parameters the attention they deserve (no more, no
less). *Journal of Vision*, *13*(7), 1–17. doi: 10.1167/13.7.3
- Roberts, G. O., Gelman, A., & Gilks, W. R. (1997). Weak conver-
gence and optimal scaling of random walk Metropolis algo-
rithms. *Annals of Applied Probability*, *7*(1), 110–120. doi:
10.1214/aoap/1034625254
- Rosenthal, J. S. (2011). Optimal proposal distributions and adap-
tive mcmc. In S. Brooks, A. Gelman, G. Jones, & X.-
L. Meng (Eds.), *Handbook of Markov Chain Monte Carlo*
(pp. 93–112). Chapman and Hall CRC. doi: 10.1201/
b10905
- Schütt, H. H., Harmeling, S., Macke, J. H., & Wichmann, F. A.
(2016). Painfree and accurate bayesian estimation of psycho-
metric functions for (potentially) overdispersed data. *Vision*
Research, *122*, 105–123. doi: 10.1016/j.visres.2016.02.002
- Scott, B. B., Constantinople, C. M., Erlich, J. C., Tank, D. W., &
Brody, C. D. (2015). Sources of noise during accumulation
of evidence in unrestrained and voluntarily head-restrained
rats. *eLife*, *4*, e11308. doi: 10.7554/eLife.11308
- Treutwein, B. (1995). Adaptive psychophysical procedures. *Vision*
Research, *35*(17), 2503–2522. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/0042-6989(95)00016-X)
0042-6989(95)00016-X

- 1176 Vul, E., Bergsma, J., & MacLeod, D. (2010). Functional adaptive
 1177 sequential testing. *Seeing and Perceiving*, 23(5), 483–515.
 1178 doi: <https://doi.org/10.1163/187847510X532694>
- 1179 Watson, A. B. (2017). QUEST+: A general multidimensional
 1180 Bayesian adaptive psychometric method. *Journal of Vision*,
 1181 17(3), 10. doi: 10.1167/17.3.10
- 1182 Watson, A. B., & Pelli, D. G. (1983). QUEST: A Bayesian adaptive
 1183 psychometric method. *Perception & psychophysics*, 33(2),
 1184 113–120. doi: 10.3758/BF03202828
- 1185 Watson, C. S., Kellogg, S. C., Kawanishi, D. T., & Lucas, P. A.
 1186 (1973). The uncertain response in detection-oriented psy-
 1187 chophysics. *Journal of Experimental Psychology*, 99(2),
 1188 180–185. doi: 10.1037/h0034736
- 1189 Wichmann, F. A., & Hill, N. J. (2001a). The psychometric
 1190 function: I. Fitting, sampling, and goodness of fit. *Per-
 1191 ception & psychophysics*, 63(8), 1293–1313. doi: 10.3758/
 1192 BF03194544
- 1193 Wichmann, F. A., & Hill, N. J. (2001b). The psychometric func-
 1194 tion: II. Bootstrap-based confidence intervals and sampling.
 1195 *Perception & Psychophysics*, 63(8), 1314–1329.
- 1196 Zocchi, S. S., & Atkinson, A. C. (1999). Optimum experimental
 1197 designs for multinomial logistic models. *Biometrics*, 55(2),
 1198 437–444. doi: 10.1111/j.0006-341X.1999.00437.x

Appendix A

Log likelihood for the classical MNL. Here we provide more
 details about the log likelihood $L = \mathbf{y}^\top \log \mathbf{p}$ under the multino-
 mial logistic model (6), first in the lapse-free case.

A convenient property of the multinomial logistic model (a prop-
 erty common to all generalized linear models) is that the parameter
 vector p_i governing y depends only on a 1-dimensional projection
 of the input, $V_i = \phi^\top \mathbf{w}_i$, which is known as the *linear predictor*.
 Recall that $\phi = \phi(\mathbf{x})$ is the input feature vector. In the multinomial
 case, it is useful to consider the column vector of linear predictors
 for a single trial, $\mathbf{V} = [V_1, \dots, V_k]^\top$, and the concatenated weight
 vector $\mathbf{w} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_k^\top]^\top$, consisting of all weights stacked
 into a single vector. We can summarize their linear relationship
 as $\mathbf{V} = X\mathbf{w}$, where X is a block diagonal matrix containing k
 blocks of ϕ^\top along the diagonal. In other words,

$$X = \begin{bmatrix} \phi^\top & \mathbf{0}^\top & \dots & \mathbf{0}^\top \\ \mathbf{0}^\top & \phi^\top & \dots & \mathbf{0}^\top \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}^\top & \mathbf{0}^\top & \dots & \phi^\top \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_k \end{bmatrix}. \quad (26)$$

Derivatives. It is convenient to work in terms of the linear pre-
 dictor $\mathbf{V} = \{V_i\}$ first. If $N_y \equiv \sum_i y_i = 1$ is the total number of
 responses per trial, the first and second derivatives of L with respect
 to \mathbf{V} are $\partial L / \partial V_j = y_j - N_y p_j$ and $\partial^2 L / \partial V_i \partial V_j = N_y p_i (\delta_{ij} - p_j)$,
 respectively. Rewriting in vector forms, we have

$$\frac{\partial L}{\partial \mathbf{V}} = (\mathbf{y} - N_y \mathbf{p})^\top, \quad (27)$$

$$\frac{\partial^2 L}{\partial \mathbf{V}^2} = -N_y (\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top) \equiv -N_y \Gamma(\mathbf{p}), \quad (28)$$

where $\text{diag}(\mathbf{p}) = [p_i \delta_{ij}]$ is a square matrix with the elements of \mathbf{p}
 on the diagonal, and zeros otherwise.

Putting back in terms of the weight vector \mathbf{w} is easy, thanks to
 the linear relationship $\mathbf{V} = X\mathbf{w}$:

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial L}{\partial \mathbf{V}} X = (\mathbf{y} - \mathbf{p})^\top X \equiv \mathbf{\Delta}^\top, \quad (29)$$

$$\frac{\partial^2 L}{\partial \mathbf{w}^2} = X^\top \frac{\partial^2 L}{\partial \mathbf{V}^2} X = -X^\top \Gamma X \equiv -\Lambda. \quad (30)$$

1230 **Concavity.** Importantly, L is concave with respect to \mathbf{V} (and
1231 therefore with respect to \mathbf{w}). To prove the concavity of L , we show
1232 that the Hessian $H = -\text{diag}(\mathbf{p}) + \mathbf{p}\mathbf{p}^\top \equiv -\Gamma$ is negative semi-
1233 definite, which is equivalent to showing that $\mathbf{z}^\top \Gamma \mathbf{z} \geq 0$:

$$\begin{aligned} \mathbf{z}^\top \Gamma \mathbf{z} &= \mathbf{z}^\top \text{diag}(\mathbf{p})\mathbf{z} - (\mathbf{z}^\top \mathbf{p})^2 \\ &= \sum_i z_i^2 p_i - \left(\sum_j z_j p_j\right)^2 \\ &= \sum_i p_i \left[\left(z_i - \sum_j z_j p_j\right)^2 \right] \geq 0 \end{aligned} \quad (31)$$

1238 for an arbitrary vector \mathbf{z} .

1239 **Log likelihood with lapse.** With a finite lapse rate λ , to recap,
1240 the multinomial logistic model is modified as $p_i = (1 - \lambda)q_i + \lambda c_i$
1241 where

$$q_i = \frac{\exp(V_i)}{\sum_j \exp(V_j)}, \quad \lambda c_i = \frac{\exp(u_i)}{1 + \sum_j \exp(u_j)}. \quad (32)$$

1243 Let us introduce the following abbreviations,

$$r_i \equiv \frac{\lambda c_i}{p_i}, \quad t_i \equiv y_i(1 - r_i), \quad s_i \equiv y_i r_i(1 - r_i), \quad (33)$$

1245 where the dimensionless ratio $r \in [0, 1]$ can be considered as the
1246 order parameter for the effect of lapse.

1247 **Derivatives with respect to the weights.** Differentiating with the
1248 linear predictor \mathbf{V} , we get

$$\begin{aligned} \frac{\partial q_i}{\partial V_l} &= (\delta_{il} - q_l)q_i, \\ \frac{\partial^2 q_i}{\partial V_j \partial V_l} &= [(\delta_{ij} - q_j)(\delta_{il} - q_l) - (\delta_{jl}q_l - q_j q_l)] q_i. \end{aligned}$$

1252 which leads to

$$\frac{\partial p_i}{\partial V_l} = (1 - \lambda) \frac{\partial q_i}{\partial V_l}, \quad \frac{\partial^2 p_i}{\partial V_j \partial V_l} = (1 - \lambda) \frac{\partial^2 q_i}{\partial V_j \partial V_l}.$$

1254 We are interested in the derivatives of the log likelihood $L =$
1255 $\mathbf{y}^\top \log \mathbf{p}$ with respect to \mathbf{V} . The partial gradient:

$$\begin{aligned} \frac{\partial L}{\partial V_l} &= \sum_i y_i \frac{1}{p_i} \frac{\partial p_i}{\partial V_l} = (1 - \lambda) \sum_i y_i \frac{q_i}{p_i} (\delta_{il} - q_l) \\ &= t_l - q_l \sum_i t_i. \end{aligned}$$

Similarly, the partial Hessian is written as

$$\begin{aligned} \frac{\partial^2 L}{\partial V_j \partial V_l} &= \sum_i y_i \left(\frac{1}{p_i} \frac{\partial^2 p_i}{\partial V_j \partial V_l} - \frac{1}{p_i^2} \frac{\partial p_i}{\partial V_j} \frac{\partial p_i}{\partial V_l} \right) \\ &= \delta_{jl} (s_l - q_l \sum_i t_i) - (q_j s_l + q_l s_j) + q_j q_l (\sum_i s_i + \sum_i t_i). \end{aligned}$$

In vector forms, and with $\tau \equiv \sum_i t_i$ and $\sigma \equiv \sum_i s_i$,

$$\frac{\partial L}{\partial \mathbf{V}} = (\mathbf{t} - \tau \mathbf{q})^\top; \quad (34)$$

$$\begin{aligned} \frac{\partial^2 L}{\partial \mathbf{V}^2} &= \text{diag}(\mathbf{s} - \tau \mathbf{q}) - (\mathbf{q}\mathbf{s}^\top + \mathbf{s}\mathbf{q}^\top) + (\tau + \sigma)\mathbf{q}\mathbf{q}^\top \\ &= -\tau [\text{diag}(\mathbf{q}) - \mathbf{q}\mathbf{q}^\top] \\ &\quad + [\text{diag}(\mathbf{s}) - (\mathbf{q}\mathbf{s}^\top + \mathbf{s}\mathbf{q}^\top) + \sigma \mathbf{q}\mathbf{q}^\top]. \end{aligned} \quad (35)$$

Note that we recover $t_i \rightarrow y_i$ and $s_i \rightarrow 0$ in the lapse-free limit
 $\lambda \rightarrow 0$. Hence the first square bracket in (35) reduces back to
the lapse-free Hessian, while the second square bracket vanishes
as $\lambda \rightarrow 0$.

In the presence of lapse, one might still be interested in the
partial Hessian with respect to the weight parameters, $H \equiv$
 $\partial^2 L / \partial \mathbf{V}^2$, which should be evaluated as in (35). To test the nega-
tive semi-definiteness of this partial Hessian, again for an arbitrary
vector \mathbf{z} , we end up with

$$\mathbf{z}^\top H \mathbf{z} = -\sum_j t_j \left\langle (z - \langle z \rangle_q)^2 \right\rangle_q + \sum_j s_j (z_j - \langle z \rangle_q)^2 \quad (36)$$

where $\langle x \rangle_q = \sum_j x_j q_j$. The partial Hessian is asymptotically neg-
ative semi-definite (which is equivalent to the log likelihood being
concave) in the lapse-free limit, where $t_j \rightarrow y_j$ and $s_j \rightarrow 0$.

Derivatives with respect to lapse parameters. From (2) and (3),
we have $p_i = (1 - \lambda)q_i + \lambda c_i$ where

$$c_i = \frac{\exp(u_i)}{\sum_j \exp(u_j)}; \quad \lambda = \frac{\sum_j \exp(u_j)}{1 + \sum_j \exp(u_j)}. \quad (37)$$

Differentiating with respect to the auxiliary lapse parameter u_i ,

$$\frac{\partial c_i}{\partial u_j} = (\delta_{ij} - c_i)c_j; \quad \frac{\partial \lambda}{\partial u_j} = (1 - \lambda)\lambda c_j. \quad (38)$$

The gradient is then

$$\frac{\partial p_i}{\partial u_j} = (\delta_{ij} - p_i)\lambda c_j; \quad (39)$$

using the abbreviations in (33), the gradient of the log likelihood is

$$\frac{\partial L}{\partial u_j} = \sum_i y_i \frac{1}{p_i} \frac{\partial p_i}{\partial u_j} = r_j (y_j - N_y \cdot p_j). \quad (40)$$

Second derivative with respect to lapse:

$$\frac{\partial^2 p_i}{\partial u_j \partial u_l} = \delta_{jl} \frac{\partial p_i}{\partial u_l} - (\delta_{ij} + \delta_{il} - 2p_i) \lambda c_l \lambda c_j; \quad (41)$$

it is useful to notice that

$$\frac{\partial p_i}{\partial u_j} \frac{\partial p_i}{\partial u_l} = \delta_{jl} \frac{\partial p_i}{\partial u_l} \lambda c_l - p_i (\delta_{ij} + \delta_{il} - 2p_i) \lambda c_l \lambda c_j. \quad (42)$$

The corresponding part of the Hessian:

$$\begin{aligned} \frac{\partial^2 L}{\partial u_j \partial u_l} &= \sum_i y_i \left(\frac{1}{p_i} \frac{\partial^2 p_i}{\partial u_j \partial u_l} - \frac{1}{p_i^2} \frac{\partial p_i}{\partial u_j} \frac{\partial p_i}{\partial u_l} \right) \\ &= \delta_{jl} \sum_i y_i \frac{1}{p_i} \left(1 - \frac{\lambda c_l}{p_i} \right) \frac{\partial p_i}{\partial u_l} \\ &= \delta_{jl} \left(s_l - r_l p_l N_y + r_l^2 p_l^2 \sum_i \frac{y_i}{p_i} \right). \end{aligned} \quad (43)$$

Finally, the mixed derivative:

$$\frac{\partial^2 p_i}{\partial u_j \partial V_l} = -(1 - \lambda) \lambda c_j \cdot (\delta_{il} - q_l) q_l. \quad (44)$$

again it is useful to notice that

$$\frac{\partial p_i}{\partial u_j} \frac{\partial p_i}{\partial V_l} = -(\delta_{ij} - p_i) \frac{\partial^2 p_i}{\partial u_j \partial V_l}. \quad (45)$$

Hence

$$\begin{aligned} \frac{\partial^2 L}{\partial u_j \partial V_l} &= \sum_i y_i \left(\frac{1}{p_i} \frac{\partial^2 p_i}{\partial u_j \partial V_l} - \frac{1}{p_i^2} \frac{\partial p_i}{\partial u_j} \frac{\partial p_i}{\partial V_l} \right) \\ &= -s_j \left(\delta_{jl} + \frac{q_l^2}{q_j} \right). \end{aligned} \quad (46)$$

From (40), (43) and (46), we see that all derivatives involving the lapse parameter scale with at least one order of r , therefore vanishing in the lapse-free limit $\lambda \rightarrow 0$.

Appendix B

The Metropolis-Hastings algorithm. The Metropolis-Hastings algorithm (Metropolis et al., 1953) generates a chain of samples, using a proposal density and a method to accept or reject the proposed moves.

A proposal is made at each iteration, where the algorithm randomly chooses a candidate for the next sample value \mathbf{x}' based on the current sample value \mathbf{x}_t . The choice follows the proposal density function, $\mathbf{x}' \sim Q(\mathbf{x}'|\mathbf{x}_t)$. When the proposal density Q is symmetric, for example a Gaussian, the sequence of samples is a random walk. In general the width of Q should match with the statistics of the distribution being sampled, and individual dimensions in the sampling space may behave differently in the multivariate case; finding the appropriate Q can be difficult.

The proposed move is either accepted or rejected with some probability; if rejected, the current sample value is reused in the next iteration, $\mathbf{x}' = \mathbf{x}_t$. The probability of acceptance is determined by comparing the values of $P(\mathbf{x}_t)$ and $P(\mathbf{x}')$, where $P(\mathbf{x})$ is the distribution being sampled. Because the algorithm only considers the acceptance ratio $\rho = P(\mathbf{x}')/P(\mathbf{x}_t) = f(\mathbf{x}')/f(\mathbf{x}_t)$ where $f(\mathbf{x})$ can be any function proportional to the desired distribution $P(\mathbf{x})$, there is no need to worry about the proper normalization of the probability distribution. If $\rho \geq 1$, the move is always accepted; if $\rho < 1$, it is accepted with a probability ρ . Consequently the samples tend to stay in the high-density regions, visiting the low-density regions only occasionally.

Optimizing the sampler. One of the major difficulties in using the MCMC method is to make an appropriate choice of the proposal distribution, which may significantly affect the performance of the sampler. If the proposal distribution is too narrow, it will take a long time for the chain to diffuse away from the starting point, producing a chain with highly correlated samples, requiring a long time to achieve independent samples. On the other hand if the proposal distribution is too wide, most of the proposed moves would be rejected, once again resulting in the chain stuck at the initial point. In either case the chain would “mix” poorly (Rosenthal, 2011). In this paper we restrict our consideration to the Metropolis-Hastings algorithm (Metropolis et al., 1953), although the issue of proposal distribution optimization is universal in most variants of MCMC algorithms, only with implementation-level differences.

The basic idea is that the optimal width of the proposal distribu-

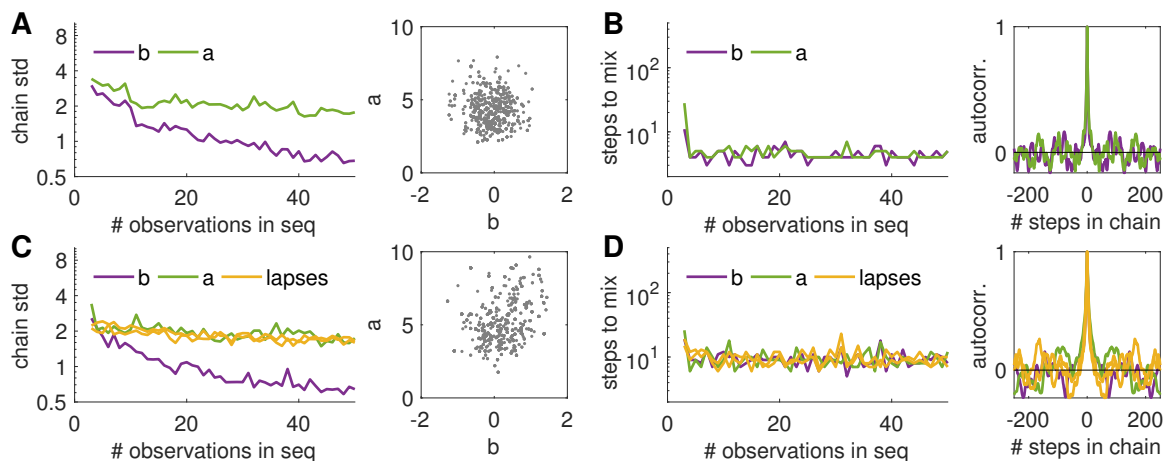


Figure 10: Statistics of the semi-adaptive MCMC in a simulated experiment, with $M = 1000$ samples per chain. We used the same binomial model as in Fig. 3, and the uniform stimulus selection algorithm. **(A-B)** In a lapse-free model: **(A)** The standard deviation of the samples, along each dimension of the parameter space, decreases as the learning progresses, as expected because the posterior distribution should narrow down as more observations are collected. Also shown is the scatter plot of all 1000 samples at the last trial $N = 50$, where the true parameter values are $(a, b) = (5, 0)$. **(B)** The mixing time of the chain (number of steps before the autocorrelation falls to $1/e$) quickly converges to some small value, meaning that the sampler is quickly optimized. Autocorrelation function at the last trial $N = 50$ is shown. **(C-D)** Same information as **(A)** and **(B)**, but with a lapse rate of $\lambda = 0.1$, with uniform lapse ($c_1 = c_2 = 1/2$).

1360 tion would be determined in proportion to the typical length scale
 1361 of the distribution being sampled. This idea was made precise in
 1362 the case of a stationary random-walk Metropolis algorithm with
 1363 Gaussian proposal distributions, by comparing the covariance ma-
 1364 trix Σ_p of the proposal distribution to the covariance matrix Σ of
 1365 the sampled chain. Once a linear scaling relation $\Sigma_p = s_d \Sigma$ is
 1366 fixed, it was observed that it is optimal to have $s_d = (2.38)^2/d$
 1367 where d is the dimensionality of the sampling space (Gelman et
 1368 al., 1996; Roberts et al., 1997). An adaptive Metropolis algo-
 1369 rithm (Haario et al., 2001) followed this observation, where the
 1370 Gaussian proposal distribution adapts continuously as the sampling
 1371 progresses. Their adaptive algorithm used the same scaling rule
 1372 $\Sigma_p = s_d \Sigma$, but updates Σ_p at each proposal where Σ is covariance
 1373 of the samples accumulated so far. Additionally, a small diagonal
 1374 component was added for stability, as $\Sigma_p = s_d(\Sigma + \epsilon I)$. We used
 1375 $\epsilon = 0.0001$ in this work.

1376 Here we propose and use the semi-adaptive Metropolis-Hastings
 1377 algorithm, which is a coarse-grained version of the original adap-
 1378 tive algorithm by Haario et al. (2001). The major difference in

our algorithm is that the adjustment of the proposal distribution is
 made only at the end of each (sequential) chain, rather than at each
 proposal within the chain. This coarse-graining is a reasonable ap-
 proximation because we will be sampling the posterior distribution
 many times as it refines over the course of data collection, once
 after each trial. Assuming that the change in posterior distribu-
 tion after each new observation is small enough, we can justify our
 use of the statistics of the previous chain to adjust the properties
 of the current chain. Unlike in the fully adaptive algorithm where
 the proposal distribution needs to stabilize quickly within a single
 chain, we can allow multiple chains until stabilization, usually a
 few initial observations – leaving some room for the coarse-grained
 approximation. This is because, for our purpose, it is not impera-
 tive that we have a good sampling of the distribution at the very
 early stage of the learning sequence where the accuracy is already
 limited by the smallness of the dataset.

When applied to the sequential learning algorithm, our semi-
 adaptive Metropolis sampler shows a consistent well-mixing prop-
 erty after a few initial adjustments, with the standard deviation

of each sampling dimension decreasing stably as data accumulate (Fig. 10). Whereas Kujala and Lukka (2006) also had the idea of adjusting the proposal density between trials, their scaling factor was fixed and independent of the sampling dimension. Building on more precise statistical observations, our method generalize well to high-dimensional parameter spaces, typical for multiple-alternative models. Our semi-adaptive sampler provides an efficient and robust alternative to the particle filter implementations (Kujala & Lukka, 2006), which has the known problem of weight degeneration (DiMattina, 2015) as the posterior distribution narrows down with the accumulation of data.

Appendix C

Fast sequential update of the posterior, with Laplace approximation. Use of Laplace approximation was shown to be particularly useful in a sequential experiment (Lewi et al., 2009), where it can be assumed that the posterior distribution after the next trial in sequence, \mathcal{P}_{t+1} , would not be very different from the current posterior \mathcal{P}_t . Let us consider the lapse-free case $\theta = \mathbf{w}$ for the moment, where the use of Laplace approximation is valid. Rearranging from (7) and (9), the sequential update for the posterior distribution is

$$\log \mathcal{P}_{t+1}(\mathbf{w}) = \log \mathcal{P}_t(\mathbf{w}) + L_{t+1}(\mathbf{w}); \quad (47)$$

or with Laplace approximation,

$$\log \mathcal{N}(\mathbf{w}|\theta_{t+1}, C_{t+1}) \approx \log \mathcal{N}(\mathbf{w}|\theta_t, C_t) + L_{t+1}(\mathbf{w}) \quad (48)$$

where $L_i(\mathbf{w}) = \log p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w})$ is a shorthand for the log likelihood of the i -th observation.

With this, we can achieve a fast sequential update of the posterior without performing the full numerical optimization each time. Because the new posterior mode θ_{t+1} is where the gradient vanishes, it can be approximated from the previous mode θ_t by taking the first derivative of (48). The posterior covariance C_{t+1} is similarly

approximated by taking the second derivate.

$$\theta_{t+1} = \theta_t + C_t \Delta_{t+1}, \quad \Delta_{t+1} = \left. \frac{\partial L_{t+1}}{\partial \mathbf{w}} \right|_{\mathbf{w}=\theta_t} \quad (49)$$

$$C_{t+1} = (C_t^{-1} + \Lambda_{t+1})^{-1}, \quad \Lambda_{t+1} = - \left. \frac{\partial^2 L_{t+1}}{\partial \mathbf{w}^2} \right|_{\mathbf{w}=\theta_{t+1}} \quad (50)$$

Using the matrix inversion lemma (Henderson & Searle, 1981), we can rewrite the posterior covariance update as

$$C_{t+1} = C_t [I - (I + \Lambda_{t+1} C_t)^{-1} \Lambda_{t+1} C_t]. \quad (51)$$

Unlike in the earlier application of this trick (Lewi et al., 2009), the covariance matrix update (50) is not a rank-one update, because of the multinomial nature of our model (our linear predictor \mathbf{y} is a vector, not a scalar as in a binary model).

Integration over the parameter space: reducing the integration space.

The evaluation of expected utility function usually involves a potentially high-dimensional integral over the parameter space. With the Gaussian approximation of the posterior, we can reduce and standardize the integration space. The process consists of three steps: diagonalization, marginalization, and standardization. First we choose a new “coordinate system” of the (say q -dimensional) weight space, such that the first k elements of the extended weight vector \mathbf{w} are coupled one-to-one to the elements of k -vector \mathbf{y} . Then we marginalize to integrate out the remaining $(q - k)$ dimensions, effectively changing the integration variable from \mathbf{w} to \mathbf{y} . Finally, we use Cholesky decomposition to standardize the normal distribution which is the posterior on \mathbf{y} . The resulting integral is still multi-dimensional, due to the multinomial nature of our model. But once the distribution is standardized, there are a number of efficient numerical integration methods that can be applied. For example, in this work, we use the Sparse Grid method (Heiss & Winschel, 2008) based on Gauss-Hermite quadrature.

Diagonalization. It is clear from (19-20) and (29-30) that all parameter-dependence in our integrand is in terms of the linear predictor $\mathbf{y} = X\mathbf{w}$. That is, we are dealing with the integral of the form

$$F = \int d\mathbf{w}' \mathcal{N}(\mathbf{w}'|\hat{\mathbf{w}}', C) \cdot f(X\mathbf{w}'), \quad (52)$$

1463 where C is the covariance matrix, and $X = \bigoplus_{j=1}^k \mathbf{g}'_j^\top$ is a fixed
 1464 matrix constructed from direct sum of k vectors. It helps to work
 1465 in a diagonalized coordinate system, so that we can separate out the
 1466 relevant dimensions of \mathbf{w} . We use the singular value decomposi-
 1467 tion of the design matrix ($X = UGV^\top$ with $U = I$ and $V = Q^\top$).
 1468 Because of the direct-sum construction, XX^\top is already diagonal,
 1469 and the left singular matrix is always I in this case. Then

$$1470 \quad G = XQ^\top = \begin{bmatrix} G_k & G_q \end{bmatrix}, \quad (53)$$

1471 where G_k is a $k \times k$ diagonal matrix and G_q is a $k \times (q - k)$
 1472 matrix of zeros. We can now denote $\mathbf{w}_k = (w_1, \dots, w_k)$ and
 1473 $\mathbf{w}_q = (w_{k+1}, \dots, w_q)$ in the diagonalized variable $\mathbf{w} = Q\mathbf{w}'$,
 1474 such that

$$1475 \quad \mathbf{w} = [\mathbf{w}_k, \mathbf{w}_q]^\top, \quad G\mathbf{w} = G_k\mathbf{w}_k = (g_1w_1, g_2w_2, \dots, g_kw_k).$$

1476 **Marginalization.** Now we have

$$1477 \quad F = \int d\mathbf{w} \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, B^{-1}) \cdot f(G\mathbf{w}), \quad B^{-1} = Q C Q^\top \quad (54)$$

1478 where B is the inverse of the *new* covariance matrix after diagonal-
 1479 ization. If we block-decompose this matrix,

$$1480 \quad B = \begin{bmatrix} B_{kk} & B_{kq} \\ B_{qk} & B_{qq} \end{bmatrix}, \quad B_{kq} = (B_{qk})^\top, \quad (55)$$

1481 the Gaussian distribution is also decomposed as

$$1482 \quad \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, B^{-1}) = \mathcal{N}(\mathbf{w}_k|\hat{\mathbf{w}}_k, B_*^{-1}) \cdot \mathcal{N}(\mathbf{w}_q|(\hat{\mathbf{w}}_q - \mathbf{b}), B_{qq}^{-1})$$

1483 where $\mathbf{b} = B_{qq}^{-1}B_{qk}\mathbf{w}_k$ and $B_* = B_{kk} - B_{kq}B_{qq}^{-1}B_{qk}$. As the
 1484 non-parallel part \mathbf{w}_q is integrated out, we have marginalized the
 1485 integral. It is useful to recall that if a variable $\mathbf{w} \sim \mathcal{N}(\hat{\mathbf{w}}, C)$ is
 1486 Gaussian distributed, its linear transform $\mathbf{y} = X\mathbf{w}$ is also Gaus-
 1487 sian distributed as $\mathbf{y} \sim \mathcal{N}(\hat{\mathbf{y}}, \Sigma)$, with $\hat{\mathbf{y}} = X\hat{\mathbf{w}}$ and $\Sigma = XCX^\top$.
 1488 Changing the integration variable to $\mathbf{y} = G_k\mathbf{w}_k$ is then straight-
 1489 forward:

$$1490 \quad F = \int d\mathbf{w}_k \mathcal{N}(\mathbf{w}_k|\hat{\mathbf{w}}_k, B_*^{-1}) \cdot f(G_k\mathbf{w}_k)
 1491 \quad = \int d\mathbf{y} \mathcal{N}(\mathbf{y}|\hat{\mathbf{y}}, \Sigma) \cdot f(\mathbf{y}), \quad \Sigma = G_k B_*^{-1} G_k^\top. \quad (56)$$

Standardization. Finally, in order to deal with the numerical in-
 1493 tegration, it is convenient to have the normal distribution standard-
 1494 ized. We can use the Cholesky decomposition for the covariance
 1495 matrix,
 1496

$$1497 \quad LL^\top = \Sigma_{t+1}, \quad (57)$$

1498 such that the new variable $\boldsymbol{\theta} = L^{-1}(\mathbf{y} - \hat{\mathbf{y}}_{t+1})$ is standard normal
 1499 distributed. From the above formulation, L can be written directly
 1500 in terms of the Cholesky decomposition of B_* :

$$1501 \quad L = G_k R^{-1} \quad \text{where} \quad R^\top R = B_*. \quad (58)$$

1502 Importantly, with this transformation, each dimension of $\boldsymbol{\theta}$ is inde-
 1503 pendently and identically distributed. The objective function to be
 1504 evaluated is now

$$1505 \quad F(\mathbf{x}) = \int d\mathbf{y} \cdot \mathcal{N}(\mathbf{y}|\hat{\mathbf{y}}_{t+1}, \Sigma_{t+1}) \cdot f(\mathbf{y}, \mathbf{x})
 1506 \quad = \int d\boldsymbol{\theta} \cdot \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, I) \cdot f(\phi(\boldsymbol{\theta}), \mathbf{x}) \quad (59)$$

1507 where $\phi(\boldsymbol{\theta}) = \hat{\mathbf{y}}_{t+1} + L\boldsymbol{\theta}$. Once the integration is standardized this
 1508 way, there are a number of efficient numerical methods that can be
 1509 applied.
 1510