Optimal inference of sameness Supporting information

Contents

1 Decision rules of the optimal observer						
	1.1	Unequal reliabilities	. 1			
	1.2	Equal reliabilities	. 5			
2	Re	sponse probabilities of the optimal observer	. 6			
	2.1	Equal reliabilities	. 6			
	2.2	Unequal reliabilities	. 7			
	2.3	Proof that A is non-negative definite in the unequal-reliabilities case	. 8			
3 Suboptimal models		poptimal models	11			
	3.1	Response probabilities for the single-criterion and blockwise-criterion models	11			
	3.2	Single-criterion model for Experiment 2	11			
	3.3	Blockwise-criterion model for Experiment 2	12			
	3.4	Maximum-absolute-difference models	12			
4	Ba	Bayesian model comparison				
5	Simulations of animal cognition experiments					
6	Me	Methods for Experiment 1, color				
7	Supplementary Tables					
8	Supplementary Figures					

1 Decision rules of the optimal observer

Since equal reliabilities (Experiment 1) are a special case of unequal reliabilities (Experiment 2), we treat the latter first.

1.1 Unequal reliabilities

The generative model of the task is given in Figure 2A. The variables are as follows: *C* is a binary variable that denotes sameness (+1 for same, -1 for different), **s** denotes the vector of *N* orientations presented, and **x** denotes the corresponding vector of *N* internal representations. Each x_i is drawn from a Gaussian distribution with mean s_i and standard deviation σ_i . The Gaussian assumption is common, and reflects the presence of many independent, unbiased sources of noise (1). Although orientation space is circular, we can treat it as the real line, because the internal representations are close to the true stimulus and σ_i is much smaller than the

circle circumference. We refer to $1/\sigma_i^2$ as the reliability of the *i*th observation. When the stimuli are the same, each orientation s_i is equal to μ . When the stimuli are different, each s_i is drawn independently from a Gaussian distribution with mean μ and standard deviation σ_s . Regardless of sameness, the value of μ is drawn from a uniform distribution on [-L,L] (one could take $L=\pi/2$ here, but any *L* works).

Optimal (Bayesian) observers base their decision on the posterior probability distribution over the variable of interest, here *C*, given single-trial observations, here $\mathbf{x}=(x_1,...,x_N)^T$. Since *C* is a binary random variable, this posterior reduces to a single number, which can be expressed in many ways. A particularly convenient way is to express it as a log posterior ratio:

$$d = \log \frac{p(C=1|\mathbf{x})}{p(C=-1|\mathbf{x})} = \log \frac{p(\mathbf{x}|C=1)}{p(\mathbf{x}|C=-1)} + \log \frac{p(C=1)}{p(C=-1)}.$$
 (S1)

Evaluating the likelihoods in this expression, $p(\mathbf{x}|C)$ for $C=\pm 1$, requires marginalization over the stimulus orientations, $\mathbf{s}=(s_1,\ldots,s_N)^T$, and their mean, μ :

$$p(\mathbf{x} | C) = \iint p(\mathbf{x} | \mathbf{s}) p(\mathbf{s} | C, \mu) p(\mu) d\mathbf{s} d\mu.$$
(S2)

We assume that the standard deviation, σ_i , of the noise associated with a stimulus is known to the observer for each stimulus and each trial, as would be the case in a Poisson-like population code (2). Therefore, we do not need to marginalize over σ_i , but can treat it as a known parameter. We now evaluate Eq. (S2) by substituting the known distributions. If **1** denotes the column vector of length *N* with all entries equal to 1, and $\delta(x)$ is the Dirac Delta distribution, then

$$p(\mathbf{x} | C = 1) = \frac{1}{2L} \iint p(\mathbf{x} | \mathbf{s}) \delta(\mathbf{s} - \mu \mathbf{1}) d\mathbf{s} d\mu$$

$$= \frac{1}{2L} \int p(\mathbf{x} | \mathbf{s} = \mu \mathbf{1}) d\mu$$

$$= \frac{1}{2L} \int \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma_i^2}\right) d\mu$$

$$= \frac{1}{2L} \int \frac{1}{\prod_{i=1}^{N} \sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2} \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{\sigma_i^2}\right) d\mu.$$
 (S3)

In the rest of this document, all sums and products run from 1 to N unless mentioned otherwise. We rewrite the sum in the exponent as

$$\sum_{i} \frac{(x_{i} - \mu)^{2}}{\sigma_{i}^{2}} = \sum_{i} \frac{x_{i}^{2}}{\sigma_{i}^{2}} - 2\mu \sum_{i} \frac{x_{i}}{\sigma_{i}^{2}} + \mu^{2} \sum_{i} \frac{1}{\sigma_{i}^{2}}$$

$$= \left(\sum_{i} \frac{1}{\sigma_{i}^{2}}\right) \left(\mu - \frac{\sum_{i} \frac{x_{i}}{\sigma_{i}^{2}}}{\sum_{i} \frac{1}{\sigma_{i}^{2}}}\right)^{2} + V_{1},$$
(S4)

with

$$V_1 = \sum_i \frac{x_i^2}{\sigma_i^2} - \frac{\left(\sum_i \frac{x_i}{\sigma_i^2}\right)^2}{\sum_i \frac{1}{\sigma_i^2}}$$

If all reliabilities are equal, so that $\sigma_i = \sigma$, then V_1 simplifies to

$$V_1 = \frac{1}{\sigma^2} \sum_i x_i^2 - \frac{1}{N\sigma^2} \left(\sum_i x_i \right)^2 = \frac{N \text{Var } \mathbf{x}}{\sigma^2}.$$

Here, Var **x** is defined as Var $\mathbf{x} = \frac{1}{N} \sum_{i} x_i^2 - \frac{1}{N^2} \left(\sum_{i} x_i \right)^2$. For the general case, in which each item comes with its own reliability, we substitute Eq. (S4) into Eq. (S3) to find

$$p(\mathbf{x} | C = 1) = \frac{\exp\left(-\frac{V_1}{2}\right)}{2L(2\pi)^{\frac{N}{2}}\prod_i \sigma_i} \int \exp\left(-\frac{1}{2}\left(\sum_i \frac{1}{\sigma_i^2}\right) \left(\mu - \frac{\sum_i \frac{x_i}{\sigma_i^2}}{\sum_i \frac{1}{\sigma_i^2}}\right)^2\right) d\mu$$
$$= \frac{\exp\left(-\frac{V_1}{2}\right)}{2L(2\pi)^{\frac{N-1}{2}}\prod_i \sigma_i} \left(\sum_i \frac{1}{\sigma_i^2}\right)^{\frac{1}{2}},$$

where in the second equality, we assumed that the domain of μ is large compared to $\left(\sum_{i} \frac{1}{\sigma_{i}^{2}}\right)^{-\frac{1}{2}}$,

and so the integral can be evaluated over the entire real line. This assumption is justified here, because in our task the maximum-likelihood estimate of a stimulus is distributed around the true value of the stimulus in a range that is much narrower than the entire circular space.

Starting from Eq. (S2), we repeat this calculation for the hypothesis that the stimuli are different, that is C = -1:

$$p(\mathbf{x} | C = -1) = \frac{1}{2L} \int_{-L}^{L} \int_{-L} \left(\prod_{i} \frac{1}{\sqrt{2\pi\sigma_{i}^{2}}} \exp\left(-\frac{(x_{i} - s_{i})^{2}}{2\sigma_{i}^{2}}\right) \right) \left(\prod_{i} \frac{1}{\sqrt{2\pi\sigma_{s}^{2}}} \exp\left(-\frac{(s_{i} - \mu)^{2}}{2\sigma_{s}^{2}}\right) \right) d\mathbf{s} d\mu$$

$$= \frac{1}{2L} \int_{-L}^{L} \prod_{i} \left(\int \frac{1}{\sqrt{2\pi\sigma_{i}^{2}}} \exp\left(-\frac{(x_{i} - s_{i})^{2}}{2\sigma_{i}^{2}}\right) \frac{1}{\sqrt{2\pi\sigma_{s}^{2}}} \exp\left(-\frac{(s_{i} - \mu)^{2}}{2\sigma_{s}^{2}}\right) ds_{i} \right) d\mu$$

$$= \frac{1}{2L} \int_{-L}^{L} \prod_{i} \left(\frac{1}{\sqrt{2\pi(\sigma_{i}^{2} + \sigma_{s}^{2})}} \exp\left(-\frac{(x_{i} - \mu)^{2}}{2(\sigma_{i}^{2} + \sigma_{s}^{2})}\right) \right) d\mu$$

$$= \frac{\exp\left(-\frac{V_{-1}}{2}\right)}{2L(2\pi)^{\frac{N}{2}} \prod_{i} \sqrt{\sigma_{i}^{2} + \sigma_{s}^{2}}} \int_{-L}^{L} \exp\left(-\frac{1}{2} \left(\sum_{i} \frac{1}{\sigma_{i}^{2} + \sigma_{s}^{2}}\right) \left(\mu - \frac{\sum_{i} \frac{x_{i}}{\sigma_{i}^{2} + \sigma_{s}^{2}}}{\sum_{i} \frac{1}{\sigma_{i}^{2} + \sigma_{s}^{2}}}\right)^{\frac{1}{2}} d\mu$$

$$= \frac{\exp\left(-\frac{V_{-1}}{2}\right)}{2L(2\pi)^{\frac{N-1}{2}} \prod_{i} \sqrt{\sigma_{i}^{2} + \sigma_{s}^{2}}} \left(\sum_{i} \frac{1}{\sigma_{i}^{2} + \sigma_{s}^{2}}\right)^{\frac{1}{2}},$$

where

$$V_{-1} = \sum_{i} \frac{x_{i}^{2}}{\sigma_{i}^{2} + \sigma_{s}^{2}} - \frac{\left(\sum_{i} \frac{x_{i}}{\sigma_{i}^{2} + \sigma_{s}^{2}}\right)^{2}}{\sum_{i} \frac{1}{\sigma_{i}^{2} + \sigma_{s}^{2}}}.$$

The log likelihood ratio, the first term on the right-hand side of Eq. (S1), now becomes

$$\log \frac{p(\mathbf{x} | C = 1)}{p(\mathbf{x} | C = -1)} = -\frac{V_1}{2} + \frac{V_{-1}}{2} + \frac{1}{2} \sum_i \log \frac{\sigma_i^2 + \sigma_s^2}{\sigma_i^2} - \frac{1}{2} \log \frac{\sum_i \frac{1}{\sigma_i^2}}{\sum_i \frac{1}{\sigma_i^2 + \sigma_s^2}}$$
$$= \frac{1}{2} \sum_i (\tilde{w}_i - w_i) x_i^2 + \frac{\left(\sum_i w_i x_i\right)^2}{2\sum_i w_i} - \frac{\left(\sum_i \tilde{w}_i x_i\right)^2}{2\sum_i \tilde{w}_i} + \frac{1}{2} \sum_i \log \frac{w_i}{\tilde{w}_i} - \frac{1}{2} \log \frac{\sum_i w_i}{\sum_i \tilde{w}_i}$$
$$= \frac{1}{2} \left(\sum_i (\tilde{w}_i - w_i) x_i^2 + \frac{\sum_i w_i w_i x_i x_j}{\sum_i w_i} - \frac{\sum_i \tilde{w}_i \tilde{w}_j x_i x_j}{\sum_i \tilde{w}_i} + \sum_i \log \frac{w_i}{\tilde{w}_i} - \log \frac{\sum_i w_i}{\sum_i \tilde{w}_i} \right),$$

where \mathbf{w} and $\tilde{\mathbf{w}}$ are column vectors with entries

$$w_i = \frac{1}{\sigma_i^2}, \quad \tilde{w}_i = \frac{1}{\sigma_i^2 + \sigma_s^2}$$

Therefore, the log posterior ratio, d, is a quadratic form in x. In matrix notation, we can write

$$d = \frac{1}{2} \left(-\mathbf{x}^{\mathrm{T}} \mathbf{A}(\mathbf{\sigma}) \mathbf{x} + \sum_{i} \log \frac{w_{i}}{\tilde{w}_{i}} - \log \frac{\sum_{i} w_{i}}{\sum_{i} \tilde{w}_{i}} \right) + \log \frac{p_{\mathrm{same}}}{1 - p_{\mathrm{same}}},$$
(S5)

where $p(C=1)=p_{\text{same}}$ and $A(\sigma)$ is a symmetric, reliability-dependent matrix with entries

$$A_{ij} = \left(w_i - \tilde{w}_i\right)\delta_{ij} + \frac{\tilde{w}_i\tilde{w}_j}{\sum_k \tilde{w}_k} - \frac{w_iw_j}{\sum_k w_k}.$$
(S6)

Maximum-a-posteriori estimation amounts to responding "same" whenever d is positive. This leads to Eq. (2) in the main text.

1.2 Equal reliabilities

When all reliabilities are equal, $\sigma_i = \sigma$ and therefore $w_i = w$ and $\tilde{w}_i = \tilde{w}$, as in Experiment 1, the entries of **A** reduce to $A_{ij} = (w - \tilde{w}) \left(\delta_{ij} - \frac{1}{N} \right)$, and the quadratic form $\mathbf{x}^T \mathbf{A} \mathbf{x}$ becomes $N(w - \tilde{w}) \operatorname{Var} \mathbf{x}$. The decision rule, which describes when the optimal observer responds "C=1", takes the form

$$N \operatorname{Var} \mathbf{x} < \frac{1}{w - \tilde{w}} \left[(N - 1) \log \frac{w}{\tilde{w}} + 2 \log \frac{p_{\text{same}}}{1 - p_{\text{same}}} \right].$$
(S7)

Thus, the decision rule amounts to comparing the variance to a criterion that depends in a specific way on the common reliability of the items (hidden in w and \tilde{w}) and on set size, N.

2 Response probabilities of the optimal observer

In the previous section, we derived the decision rules of the optimal observers for both experiments. However, these decision rules cannot be tested directly, since as experimenters we do not have access to the internal observations \mathbf{x} . We only know the presented stimuli, \mathbf{s} , and the subject's response on each trial. To obtain testable model predictions, we compute the theoretical probability of a "same" response given by the optimal observer, when the observations \mathbf{x} are drawn from the generative model for given \mathbf{s} . In the resulting model predictions, \mathbf{x} is no longer present. The complete model thus consists of two parts "merged together": the generative model, producing a theoretical (or simulated) set of observations, and the inference model, which describes the decision made by the optimal observer based on those observations.

Formally, we denote the observer's estimate of C by \hat{C} ; for the optimal observer, it equals the sign of d. The optimal observer responds according to this estimate. Eq. (S5) determines the optimal decision rule on a single trial. We denote the probability of the optimal observer responding \hat{C} when the experimenter presents a specific stimulus set \mathbf{s} by $p(\hat{C}|\mathbf{s})$. It is obtained by averaging over the hypothetical observations \mathbf{x} generated by \mathbf{s} :

$$p(\hat{C} | \mathbf{s}) = \int p(\hat{C} | \mathbf{x}) p(\mathbf{x} | \mathbf{s}) d\mathbf{x} = \int \delta_{\hat{C}, \operatorname{sgn}(d(\mathbf{x}))} p(\mathbf{x} | \mathbf{s}, \mathbf{\sigma}) d\mathbf{x}, \qquad (S8)$$

where δ is the Kronecker delta. In words, the probability of responding $\hat{C} = 1$ is equal to the proportion of observations drawn from $p(\mathbf{x}|\mathbf{s})$ for which $d(\mathbf{x})$ is positive. For many psychophysical tasks, an integral like that given in Eq. (S8) is analytically intractable, and needs to be approximated, for example, by Monte Carlo simulation. In the present task, an analytical treatment is possible.

2.1 Equal reliabilities

We first consider the case of Experiment 1, with the decision rule given by Eq. (S7). We are interested in the probability that this inequality is satisfied for **x** drawn from a multivariate normal distribution with mean **s** and covariance matrix $\sigma^2 \mathbf{I}$, denoted $\mathbf{x} \sim \mathcal{N}(\mathbf{s}, \sigma^2 \mathbf{I})$, with **I** the identity matrix. We write Var $\mathbf{x}=\mathbf{x}^T \mathbf{A}'\mathbf{x}$, where the components of \mathbf{A}' are $A'_{ij}=\delta_{ij}/N-1/N^2$. This matrix has N-1 eigenvalues equal to 1, and one eigenvalue equal to 0. Therefore, for an orthogonal matrix **M** and diagonal matrix $\mathbf{D}=\text{diag}(1,1,\ldots,1,0)$, we can write $\mathbf{M}^T\mathbf{A}'\mathbf{M}=\mathbf{D}$. Let **Y** be $\mathbf{Y}=\mathbf{M}^T\mathbf{x}$. Since **M** is orthogonal, $\mathbf{y}\sim \mathcal{N}(\mathbf{\mu},\sigma^2 \mathbf{I})$, with $\mathbf{\mu}=\mathbf{M}^T\mathbf{s}$. Since

$$\sum_{i=1}^{N-1} \frac{y_i^2}{\sigma^2} = \frac{1}{\sigma^2} \mathbf{Y}^{\mathsf{T}} \mathbf{D} \mathbf{Y} = \frac{1}{\sigma^2} \mathbf{x}^{\mathsf{T}} \mathbf{M} \mathbf{D} \mathbf{M}^{\mathsf{T}} \mathbf{x} = \frac{1}{\sigma^2} \mathbf{x}^{\mathsf{T}} \mathbf{A}' \mathbf{x} , \qquad (S9)$$

it follows that $\mathbf{x}^{\mathrm{T}}\mathbf{A}'\mathbf{x}/\sigma^2$ follows a non-central chi-squared distribution with N-1 degrees of

freedom. This variable is equal to *wN*Var **x**. The noncentrality parameter of the sum on the left of Eq. (S9) is $\sum_{i=1}^{N-1} \frac{\mu_i^2}{\sigma^2}$, which can be rewritten as

$$\sum_{i=1}^{N-1} \frac{\mu_i^2}{\sigma^2} = \frac{1}{\sigma^2} \boldsymbol{\mu}^{\mathrm{T}} \mathbf{D} \boldsymbol{\mu} = \frac{1}{\sigma^2} \mathbf{s}^{\mathrm{T}} \mathbf{M} \mathbf{D} \mathbf{M}^{\mathrm{T}} \mathbf{s} = \frac{1}{\sigma^2} \mathbf{s}^{\mathrm{T}} \mathbf{A}' \mathbf{s} = wN \operatorname{Var} \mathbf{s}$$

Therefore, we have

$$wN \operatorname{Var} \mathbf{x} \sim \chi^2_{N-1} (wN \operatorname{Var} \mathbf{s}).$$
(S10)

From Eq. (S7), the probability of responding "same" equals the probability that

wN Var **x** is less than $\frac{w}{w - \tilde{w}} \left[(N-1) \log \frac{w}{\tilde{w}} + 2 \log \frac{p_{\text{same}}}{1 - p_{\text{same}}} \right]$. Using Eq. (S10), this probability is

obtained from a cumulative chi-squared distribution:

$$p(\hat{C}=1|\mathbf{s}) = \Pr\left(\chi_{N-1}^{2}(wN\operatorname{Var}\mathbf{s}) < \frac{w}{w-\tilde{w}}\left[(N-1)\log\frac{w}{\tilde{w}} + 2\log\frac{p_{\text{same}}}{1-p_{\text{same}}}\right]\right).$$
(S11)

2.2 Unequal reliabilities

We next consider the case of Experiment 2, where reliability can differ between stimuli. Then **x** follows a multivariate normal distribution with mean **s** and covariance matrix $\Sigma = \text{diag}(\sigma_1^2, ..., \sigma_N^2)$. We define a random variable $Q(\mathbf{x}) = \mathbf{x}^T \mathbf{B} \mathbf{x}$, where **B** is a non-negative definite and symmetric matrix. We denote mean and standard deviation of Q by μ_Q and σ_Q , respectively. Liu, Tang, and Zhang (3) have found the following approximation to the distribution of Q:

$$\Pr\left(\mathbf{x}^{\mathrm{T}}\mathbf{B}\mathbf{x} < k\right) \approx \Pr\left(\chi_{l}^{2}\left(\delta\right) < t\sigma_{\chi} + \mu_{\chi}\right)$$
(S12)

where the variables *l*, δ , *t*, σ_{χ} , and μ_{χ} are defined as follows:

$$t = \frac{k - \mu_{Q}}{\sigma_{Q}}$$

$$l = 3a^{2} - 2r_{1}a^{3}$$

$$\delta = r_{1}a^{3} - a^{2}$$

$$\mu_{\chi} = 2a^{2} - r_{1}a^{3}$$

$$\sigma_{\chi} = a\sqrt{2}$$

$$a = \begin{cases} \frac{1}{r_{1} - \sqrt{r_{1}^{2} - r_{2}}}, & \text{if } r_{1}^{2} > r_{2} \\ \frac{1}{r_{1}}, & \text{if } r_{1}^{2} \le r_{2} \end{cases}$$

$$r_{1} = \frac{c_{3}}{c_{2}^{3/2}}$$

$$r_{2} = \frac{c_{4}}{c_{2}^{2}}$$

$$c_{j} = \operatorname{Tr}((\mathbf{B}\mathbf{\Sigma})^{j}) + j\mathbf{s}^{\mathrm{T}}(\mathbf{B}\mathbf{\Sigma})^{j-1}\mathbf{B}\mathbf{s}$$

Numerical simulation shows that for our purposes, this is a very good approximation. In order to apply this approximation to $\mathbf{x}^{T}\mathbf{A}\mathbf{x}$ from Eq.(S5), **A** has to be both symmetric and non-negative definite. Since **A** is obviously symmetric, what remains to be shown is that in our case, **A** is non-negative definite.

2.3 **Proof that A is non-negative definite in the unequal-reliabilities case**

The matrix A specified by Eq. (S6) can be written as

$$\mathbf{A} = \operatorname{diag}(\mathbf{w} - \tilde{\mathbf{w}}) + \frac{\tilde{\mathbf{w}}\tilde{\mathbf{w}}^{\mathrm{T}}}{\tilde{\mathbf{w}}^{\mathrm{T}}\mathbf{1}} - \frac{\mathbf{w}\mathbf{w}^{\mathrm{T}}}{\mathbf{w}^{\mathrm{T}}\mathbf{1}}$$

To show that *A* is non-negative definite (positive semidefinite), we prove that its eigenvalues are non-negative. In Experiment 2, *m* of the σ_i 's are equal to σ_{low} , and *N*-*m* are equal to σ_{high} . Without loss of generality, we arrange them as $\sigma = (\sigma_{low}, ..., \sigma_{low}, \sigma_{high}, ..., \sigma_{high})$, *i.e.*, the first *m* correspond to low-reliability stimuli, and the last *N*-*m* to high-reliability ones. Consequently, $w_1 = ... = w_m = w_{low} = 1/\sigma_{low}^2$ and $w_{m+1} = ... = w_N = w_{high} = 1/\sigma_{high}^2$. The $(i,j)^{\text{th}}$ term of **A** is given by Eq. (S6). Note that

$$\left(\mathbf{A1}\right)_{i} = \sum_{j=1}^{N} A_{ij} = w_{i} - \tilde{w}_{i} + \sum_{j=1}^{N} \left(\frac{\tilde{w}_{i} \tilde{w}_{j}}{\sum_{k} \tilde{w}_{k}} - \frac{w_{i} w_{j}}{\sum_{k} w_{k}} \right) = 0, \qquad (S13)$$

we see that $\mathbf{1}=(1,...,1)^{T}$ is an eigenvector with eigenvalue 0. We now define a set of m-1 vectors \mathbf{v}^{k} , with k=2,...,m, as follows:

$$v_i^k = \delta_{i1} - \delta_{ik}$$
 .

For example, $\mathbf{v}^3 = (1,0,-1,0,...0)^T$. These vectors are eigenvectors of **A** with eigenvalues $w_{\text{low}} - \tilde{w}_{\text{low}}$, since

$$\begin{split} \left(\mathbf{A}\mathbf{v}^{k}\right)_{i} &= \sum_{j=1}^{N} A_{ij} v_{j}^{k} = \sum_{j=1}^{N} A_{ij} \left(\delta_{j1} - \delta_{jk}\right) \\ &= \sum_{j=1}^{N} \left(\left(w_{i} - \tilde{w}_{i}\right) \delta_{ij} + \frac{\tilde{w}_{i} \tilde{w}_{j}}{\sum_{k} \tilde{w}_{k}} - \frac{w_{i} w_{j}}{\sum_{k} w_{k}} \right) \left(\delta_{j1} - \delta_{jk}\right) \\ &= \sum_{j=1}^{N} \left(\left(w_{i} - \tilde{w}_{i}\right) \delta_{ij} \left(\delta_{j1} - \delta_{jk}\right) + \left(\frac{\tilde{w}_{i} \tilde{w}_{j}}{\sum_{k} \tilde{w}_{k}} - \frac{w_{i} w_{j}}{\sum_{k} w_{k}}\right) \left(\delta_{j1} - \delta_{jk}\right) \right) \\ &= \left(w_{i} - \tilde{w}_{i}\right) \left(\delta_{i1} - \delta_{ik}\right) + \frac{\tilde{w}_{i} \left(\tilde{w}_{1} - \tilde{w}_{k}\right)}{\sum_{k} \tilde{w}_{k}} - \frac{w_{i} \left(w_{1} - w_{k}\right)}{\sum_{k} w_{k}} \\ &= \left(w_{\mathrm{low}} - \tilde{w}_{\mathrm{low}}\right) \left(\delta_{i1} - \delta_{ik}\right) = \left(w_{\mathrm{low}} - \tilde{w}_{\mathrm{low}}\right) v_{i}^{k}. \end{split}$$

Similarly, the *N*-*m*-1 vectors \mathbf{v}^k with k=m+2,...,N, whose i^{th} entries are $v_i^k = \delta_{i,m+1} - \delta_{ik}$ are eigenvectors with eigenvalues $w_{\text{high}} - \tilde{w}_{\text{high}}$. That leaves one eigenvalue to find. Consider the vector $\mathbf{v}=(N-m,...,N-m,-m,...,-m)^{\text{T}}$. We will show that this vector is an eigenvector. We first calculate the product of any of the first *m* rows of *A* with \mathbf{v} , making use of Eq. (S13)

$$\begin{split} \left(\mathbf{Av}\right)_{i} &= \left(N-m\right) \sum_{j=1}^{m} A_{ij} - m \sum_{j=m+1}^{N} A_{ij} \\ &= \left(N-m\right) \left(w_{\text{low}} - \tilde{w}_{\text{low}}\right) + m \left(N-m\right) \left(\frac{\tilde{w}_{\text{low}}^{2}}{\sum_{k} \tilde{w}_{k}} - \frac{w_{\text{low}}^{2}}{\sum_{k} w_{k}}\right) - m \left(N-m\right) \left(\frac{\tilde{w}_{\text{low}} \tilde{w}_{\text{high}}}{\sum_{k} \tilde{w}_{k}} - \frac{w_{\text{low}} w_{\text{high}}}{\sum_{k} \tilde{w}_{k}}\right) \\ &= -\left(N-m\right) \left(N-m\right) \left(\frac{\tilde{w}_{\text{low}} \tilde{w}_{\text{high}}}{\sum_{k} \tilde{w}_{k}} - \frac{w_{\text{low}} w_{\text{high}}}{\sum_{k} w_{k}}\right) - m \left(N-m\right) \left(\frac{\tilde{w}_{\text{low}} \tilde{w}_{\text{high}}}{\sum_{k} \tilde{w}_{k}}\right) \\ &= N \left(N-m\right) \left(\frac{w_{\text{low}} w_{\text{high}}}{\sum_{k} w_{k}} - \frac{\tilde{w}_{\text{low}} \tilde{w}_{\text{high}}}{\sum_{k} \tilde{w}_{k}}\right) \\ &= N \left(N-m\right) \left(\frac{1}{m\sigma_{\text{low}}^{2} + \left(N-m\right)\sigma_{\text{high}}^{2}} - \frac{1}{m \left(\sigma_{\text{low}}^{2} + \sigma_{s}^{2}\right) + \left(N-m\right) \left(\sigma_{\text{high}}^{2} + \sigma_{s}^{2}\right)}\right) \\ &= \left(N-m\right) \left(\frac{N^{2}\sigma_{s}^{2}}{\left(m\sigma_{\text{low}}^{2} + \left(N-m\right)\sigma_{\text{high}}^{2}\right) \left(m\sigma_{\text{low}}^{2} + \left(N-m\right)\sigma_{\text{high}}^{2} + N\sigma_{s}^{2}\right)}\right). \end{split}$$

Similarly, for any of the last *N*–*m* rows of *A*, we get

$$\left(\mathbf{Av}\right)_{i} = -m \left(\frac{N^{2}\sigma_{s}^{2}}{\left(m\sigma_{low}^{2} + \left(N - m\right)\sigma_{high}^{2}\right)\left(m\sigma_{low}^{2} + \left(N - m\right)\sigma_{high}^{2} + N\sigma_{s}^{2}\right)}\right).$$

Thus, $\frac{N^2 \sigma_s^2}{\left(m\sigma_{\text{low}}^2 + (N-m)\sigma_{\text{high}}^2\right)\left(m\sigma_{\text{low}}^2 + (N-m)\sigma_{\text{high}}^2 + N\sigma_s^2\right)}$ is the final eigenvalue of A. Since

all eigenvalues are non-negative, A is non-negative definite, and therefore we may apply the approximation discussed in Section 2.2.

3 Suboptimal models

3.1 Response probabilities for the single-criterion and blockwise-criterion models

Calculating a model's predictions for the probability of responding "same" given a set of presented orientations s consists of two steps: to determine the decision rule for responding "same", and to apply this decision rule to the set of internal representations on each individual trial. In non-optimal observer models, the first step (decision rule) is different, but the second step (calculating the probability with which the decision is satisfied) is identical. In particular, as long as the decision rule is of the form

$$\mathbf{x}^{\mathrm{T}}\mathbf{B}\mathbf{x} < k \,, \tag{S14}$$

with **B** a non-negative definite, symmetric matrix, we can still use Eq. (S12) to approximate the model response probabilities. The only difference from the optimal model is that different expressions must be substituted for **B** and/or *k*. In the special case that the decision rule is Var $\mathbf{x} < k$ and the reliabilities are equal, the model response probabilities are given by an exact expression analogous to Eq. (S11),

$$p(\hat{C}=1 | \mathbf{s}) = \Pr(\chi^2_{N-1}(wN \operatorname{Var} \mathbf{s}) < wk).$$

3.2 Single-criterion model for Experiment 2

In Experiment 2, the SC observer assumes in the decision rule that all reliabilities are equal, $\sigma_i = \sigma_{assumed}$ for all *i*. In other words, this observer does not weight the observations by their correct respective reliabilities. We will show here that this model is equivalent to one in which the observer compares the sample variance to a single criterion.

The decision rule is the same one as in Experiment 1, Eq. (S7), but with $\sigma_{assumed}$ instead of σ . This rule can be rewritten as

$$\operatorname{Var} \mathbf{x} < \frac{\sigma_{\operatorname{assumed}}^2}{N} \left(1 + \frac{\sigma_{\operatorname{assumed}}^2}{\sigma_{\operatorname{s}}^2} \right) \left[(N-1) \log \left(1 + \frac{\sigma_{\operatorname{s}}^2}{\sigma_{\operatorname{assumed}}^2} \right) + 2 \log \frac{p_{\operatorname{same}}}{1 - p_{\operatorname{same}}} \right].$$
(S15)

Since $\sigma_{assumed}$ is a free parameter, the right-hand side can assume any value on the real line. To see this, first note that this expression is continuous. Then compute the two limits

$$\lim_{\sigma_{assumed}^{2} \to \infty} \frac{\sigma_{assumed}^{2}}{N} \left(1 + \frac{\sigma_{assumed}^{2}}{\sigma_{s}^{2}} \right) \left[(N-1) \log \left(1 + \frac{\sigma_{s}^{2}}{\sigma_{assumed}^{2}} \right) + 2 \log \frac{p_{same}}{1 - p_{same}} \right] = \infty$$

$$\lim_{\sigma_{assumed}^{2} \to 0} \frac{\sigma_{assumed}^{2}}{N} \left(1 + \frac{\sigma_{assumed}^{2}}{\sigma_{s}^{2}} \right) \left[(N-1) \log \left(1 + \frac{\sigma_{s}^{2}}{\sigma_{assumed}^{2}} \right) + 2 \log \frac{p_{same}}{1 - p_{same}} \right]$$

$$= \frac{N-1}{N} \lim_{\sigma_{assumed}^{2} \to 0} \sigma_{assumed}^{2} \log \left(1 + \frac{\sigma_{s}^{2}}{\sigma_{assumed}^{2}} \right) = 0$$

Therefore, we can simply replace the entire right-hand side by a single free parameter *k*, i.e., Var $\mathbf{x} < k$, thereby justifying the terminology "single-criterion model". (Note that in order to obtain model predictions, we still need to fit the parameters σ_{low} and σ_{high} .) Thus, the SC model has only two variants: with and without lapse rate.

3.3 Blockwise-criterion model for Experiment 2

In Experiment 2, the BC model is the model in which $\sigma_{assumed}$ may vary by block type (LOW, MIXED, or HIGH). Then, the decision rule is equivalent to Var $\mathbf{x} < k_{block}$, where each block type has its own free parameter k_{block} . This model is analogous to the BC model in Experiment 1, with reliability condition (LOW, MIXED, or HIGH) taking the place of set size. As in the SC model, there are only two variants: with and without lapse rate.

3.4 Maximum-absolute-difference models

For the MAD models, the response probabilities cannot be calculated analytically. We estimated these probabilities numerically, by simulating 10,000 sets of *N* internal representations for each of the 2700 experimental trials, and applying the model's decision rule to each set; the frequency of "same" responses was an estimate of the probability of responding "same" on the given trial.

4 Bayesian model comparison

We denote by $p(\hat{C}_i|\mathbf{s}_i, M, \boldsymbol{\theta})$ the probability predicted by model M with parameters $\boldsymbol{\theta}$ of the subject's *actual* response on the *i*th trial, \hat{C}_i , when the presented stimuli are \mathbf{s}_i . We computed the log probability of the data given M by marginalizing over $\boldsymbol{\theta}$. For the prior over parameters, we assumed a uniform distribution, $p(\boldsymbol{\theta})=1/\text{Vol}_{\boldsymbol{\theta}}$, where $\text{Vol}_{\boldsymbol{\theta}}$ is the volume of parameter space. We calculated the parameter likelihood by assuming that the data are conditionally independent across trials: $p(\text{data}|M,\boldsymbol{\theta})=\prod_i p(\hat{C}_i|\mathbf{s}_i,M,\boldsymbol{\theta})$. For the MAD models, we approximated the marginalization through a Riemann sum; the parameter ranges were 1 to 40 for the decision criterion, 0.1° to 20.1° for σ , and 0 to 0.5 for the lapse rate, each in 35 steps. For the other models, since analytical expressions were available, we used the Laplace approximation (4); the sizes of the parameter ranges were 20° for σ , 20° for the assumed σ_s , 0.4 for p_{same} , and 0.4 for the lapse rate.

5 Simulations of animal cognition experiments

To examine whether the proportion of "different" responses from the optimal observer correlates with the entropy of the stimulus set, we simulated the stimulus sets used by Young et al. (5) (see Table 2 in their paper) and used these as input to the optimal-observer model. These sets always contained 16 items, with several subsets of identical items (e.g., 4 subsets of 4 identical items each). In our simulations, for each subset, a random stimulus value was drawn from a Gaussian distribution (with $\sigma_s=5$) and assigned to all items in that subset. A total of 1000 trials were simulated per stimulus set. On each trial, stimulus observations were simulated by adding Gaussian noise to the stimulus orientations ($\sigma=10$) and a response was generated by applying the decision rule from the optimal-observer model. Entropy and "scaled logit of percent different responses" were computed as described by Young et al. (5). To examine the effect of set size on the proportion of different responses in the optimal-observer model, we simulated 1000 trials per set size, with $\sigma_s=8$ and $\sigma=8$. To examine the effect of stimulus visibility, we varied internal noise, σ , and simulated 1000 trials per noise level (with $\sigma_s=5$, N=8, $p_{same}=0.6$, and a guessing rate of 0.25). While model parameters were separately chosen for each experiment, the observed trends were robust under a wide range of parameters values.

6 Experiment 1 (color)

Methods

The methods for Experiment 1, color, were identical to the Experiment 1 for orientation, except for the following differences. Stimuli consisted of a set of colored discs with a radius of 0.4 deg. The colors of the discs were drawn independently from 180 color values uniformly distributed along a circle of radius 50 in CIE 1976 (L^*,a^*,b^*) color space. This circle had constant luminance ($L^*=58$) and was centered at the point ($a^*=12$, $b^*=13$). The stimuli were presented on a grey background of luminance 10 cd/m². Set size was chosen randomly on each trial. In the models, "blockwise criterion" was now defined as one criterion per set size, rather than per block. The experiment consisted of 3 sessions with 6 blocks of 150 trials. Two authors and five paid, naive subjects participated in the experiment.

Results

The SC models can be ruled out but not the BC models (Fig. S1A-B; log likelihood differences between the optimal model and the SC, BC, MAD-SC, and MAD-BC models were 8.7 ± 4.9 , -3.9 ± 2.5 , 27.3 ± 6.3 , and 0.8 ± 2.3 , respectively). Even though the BC model just edges out the optimal model, the decision criteria and noise levels from the best BC model are very close to those predicted by the best optimal-model variant (Fig. S1C; *p*>0.1 for all six *t*-tests). Finally, the noise level exhibits a weak dependence on set size (Fig. S1C; power law fit yields a power of 0.23 ± 0.04).

7 Supplementary Tables

Table S1. Overview of model parameters. The models have the same parameters in both experiments, except that in Experiment 1, $\sigma = (\sigma_{N=2}, \sigma_{N=4}, \sigma_{N=8})$ and $\mathbf{k} = (k_{N=2}, k_{N=4}, k_{N=8})$, while in Experiment 2, $\sigma = (\sigma_{\text{low}}, \sigma_{\text{high}})$ and $\mathbf{k} = (k_{\text{LOW}}, k_{\text{MIXED}}, k_{\text{HIGH}})$.

Model	Free parameters
Optimal	σ
Optimal	σ, λ
Optimal	$\sigma, \sigma_{\rm s}$
Optimal	$\sigma, \sigma_{\rm s}, \lambda$
Optimal	σ, p_{same}
Optimal	σ , p_{same} , λ
Optimal	$\sigma, p_{same}, \sigma_s$
Optimal	$\boldsymbol{\sigma}, p_{\mathrm{same}}, \sigma_{\mathrm{s}}, \lambda$
Single-criterion	σ , <i>k</i>
Single-criterion	σ, k, λ
Blockwise-criterion	σ, k
Blockwise-criterion	$\sigma, \mathbf{k}, \lambda$
MAD with single criterion	σ , <i>k</i>
MAD with single criterion	σ , <i>k</i> , λ
MAD with blockwise criterion	σ, k
MAD with blockwise criterion	$\sigma, \mathbf{k}, \lambda$

Table S2. Overview of maximum-likelihood parameter values of the best fitting optimal models in Experiment 1. An empty cell indicates that the respective parameter was not a free parameter in the best-fitting model variant.

	p_{same}	$\sigma_{ m s}$	λ
AK	0.58	6.55	0.06
DS	0.42		0.04
HB			
MH			0.12
ML	0.40	11.5	
RB		11.3	
RC			0.06
TR	0.46		0.11

Table S3. Overview of maximum-likelihood parameter values of the best-fitting optimal models in Experiment 2. An empty cell indicates that the respective parameter was not a free parameter in the best-fitting model variant.

	p_{same}	$\sigma_{ m s}$	λ
BN	0.53		0.08
DB			0.11
DS		19.5	0.22
HB	0.45		0.11
KJ	0.45		0.07
MH	0.44		0.11
ML	0.35	13.7	
MV	0.46		0.04
RB	0.54	7.73	
RC	0.55		



Fig. S1 Comparison of models in Experiment 3. Circles and error bars represent mean and s.e.m. of subject data. Shaded areas represent s.e.m. of model fits. (*A*) Proportion "different" responses as a function of sample standard deviation for optimal and suboptimal models. (*B*) Bayesian model comparison. Each bar represent the log likelihood of the optimal model minus that of a suboptimal model. (*C*) The decision criteria (left) and the internal noise levels (right) for the best-fitting BC model are nearly identical to those of the best-fitting optimal-observer model. This suggests that the human criteria are close to optimal.



Fig. S2. A comparison of the decision strategies of the optimal and BC models in Experiment 2. (*A*) Distributions of the decision variable of the optimal model. As the number of high-reliability stimuli in the display, N_{high} , increases, the distributions become more separable, making the task easier. The optimal observer sets his criterion (dashed line) in such a way that performance is maximized. (*B*) Distributions of the decision variable of the BC model (Var **x**). This model uses the same criterion regardless of N_{high} (dashed line). Hence, it cannot maximize performance for all values of N_{high} . The same holds for the SC, MAD-SC, and MAD-BC models. The distributions differ between the models because the decision variables do.

- 1. Green DM & Swets JA (1966) *Signal detection theory and psychophysics* (John Wiley & Sons, Los Altos, CA).
- 2. Ma WJ, Beck JM, Latham PE, & Pouget A (2006) Bayesian inference with probabilistic population codes. *Nat Neurosci* 9(11):1432-1438.
- 3. Liu H, Tang Y, & Zhang HH (2009) A new chi-square approximation to the distribution of non-negative definite quadratic forms in non-central normal variables. *Comp Statistics and Data Analysis* 53:853-856.
- 4. MacKay DJ (2003) *Information theory, inference, and learning algorithms* (Cambridge University Press, Cambridge, UK).
- 5. Young ME & Wasserman EA (1997) Entropy detection by pigeons: response to mixed visual displays after same-different discrimination training. *J Exp Psychol Anim Behav Process* 23(2):157-170.