

The psychometric function II: bootstrap based confidence intervals and sampling

FELIX A. WICHMANN and N. JEREMY HILL

*Sensory Research Unit,
Department of Experimental Psychology,
The University of Oxford,
Oxford OX1 3UD, UK.*

Running title: The Psychometric Functions II
Keywords: Parametric Bootstrap, Monte Carlo simulations, psychometric function, modelling data, variance estimation.

Number of pages: 37
Number of figures: 7
Number of equations: 4
Number of tables: 2

Present address and address for correspondence:

Dr. Felix A. Wichmann
Department of Psychology
The University of Leuven
Tiensestraat 102
3000 Leuven, Belgium.

phone: +32-16-32.60.15
fax: +32-16-32.60.99
email: felix.wichmann@psy.kuleuven.ac.be

Abstract

The psychometric function relates an observer's performance to an independent variable, usually a physical quantity of an experimental stimulus. Even if a model is successfully fit to the data and its goodness-of-fit is acceptable, experimenters require an estimate of variability of the parameters to assess whether differences across conditions are significant. Accurate estimates of variability are difficult to obtain, however, given the typically small size of psychophysical datasets: traditional statistical techniques are only asymptotically correct and can be shown to be unreliable in some common situations. Here and in our companion paper (Wichmann & Hill, 2000) we suggest alternative statistical techniques based on Monte Carlo resampling methods. The current paper's principal topic is the estimation of the variability of fitted parameters and derived quantities such as thresholds and slopes: first, we outline the basic bootstrap procedure and argue in favour of the parametric as opposed to non-parametric bootstrap. Second, we describe how the bootstrap bridging assumption, on which the validity of the procedure depends, can be tested. Third, we show how one's choice of sampling scheme (the placement of sample points on the stimulus axis) strongly affects the reliability of bootstrap confidence intervals and we make recommendations on how to sample the psychometric function efficiently. Fourth, we show that, under certain circumstances, the (arbitrary) choice of the distribution function chosen can exert an unwanted influence on the size of the bootstrap confidence intervals obtained, and we make recommendations on how to avoid this influence. Finally, we introduce improved confidence intervals (bias corrected and accelerated) which improve on the parametric and percentile-based bootstrap confidence intervals previously used. Software implementing our methods is available.

Outline

The performance of an observer on a psychophysical task is typically summarized by reporting one or more *response thresholds*—stimulus intensities required to produce a given level of performance—and by a characterization of the rate at which performance improves with increasing stimulus intensity. These measures are derived from a *psychometric function*, which describes the dependence of an observer's performance on some physical aspect of the stimulus.

Fitting psychometric functions is a variant of the more general problem of modelling data. Modelling data is a three-step process: First, a model is chosen and the parameters are adjusted to minimize the appropriate error-metric or loss function. Second, error estimates of the parameters are derived and, third, the goodness-of-fit between model and data is assessed. This paper is concerned with the second of these steps, the estimation of variability in fitted parameters and in quantities derived from them. Our companion paper (Wichmann & Hill, 2000) illustrates how to fit psychometric functions avoiding bias resulting from stimulus-independent lapses and how to evaluate goodness-of-fit between model and data.

We advocate the use of Efron's *bootstrap method*, a particular kind of Monte Carlo technique to the problem of estimating the variability of parameters, thresholds and slopes of psychometric functions (Efron, 1979; 1982; Efron & Gong, 1983; Efron & Tibshirani, 1991; 1993). Bootstrap techniques are not without their own assumptions and potential pitfalls: in the course of this paper we shall discuss these, and examine their effect on the estimates of variability we obtain. We describe and examine the use of parametric bootstrap techniques in finding confidence intervals for thresholds and slopes. We then explore the sensitivity of the estimated confidence interval widths to (a) sampling schemes, (b) mismatch of the objective function and (c) accuracy of the originally fitted parameters. The last of these is particularly important since it provides a test of the validity of the *bridging assumption* on which the use of parametric bootstrap techniques rely.

Finally we recommend, based on the theoretical work of others, the use of a technique called *bias–correction with acceleration* (BC_a) to obtain stable and accurate confidence interval estimates.

Background

The psychometric function

Our notation will follow the conventions we have outlined in Wichmann and Hill (2000). A brief summary of terms follows.

Performance on K blocks of a constant-stimuli psychophysical experiment can be expressed using three vectors, each of length K . \mathbf{x} denotes the stimulus values used, \mathbf{n} denotes the numbers of trials performed at each point, and \mathbf{y} denotes the proportion of correct responses (in n -AFC experiments) or positive responses (single-interval or "yes/no" experiments) on each block. We often use N to refer to the total number of trials in the set, $N = \sum n_i$.

The number of correct responses $y_i n_i$ in a given block i is assumed to be the sum of random samples from a Bernoulli process with probability of success p_i . A psychometric function $\psi(x)$ is the function that relates the stimulus dimension x to the expected performance value p .

A common general form for the psychometric function is:

$$\psi(x; \alpha, \beta, \gamma, \lambda) = \gamma + (1 - \gamma - \lambda)F(x; \alpha, \beta) . \quad (1)$$

The shape of the curve is determined by our choice of a functional form for F , and by the four parameters $\{\alpha, \beta, \gamma, \lambda\}$, to which we shall refer collectively using the parameter vector $\boldsymbol{\theta}$. F is typically a sigmoidal function such as the Weibull, cumulative Gaussian, logistic or Gumbel. We assume that F describes the underlying psychological mechanism of interest: the parameters γ and λ determine the lower and upper bounds of the curve, which are affected by other factors. In yes/no paradigms γ is the "guess rate" and λ the

"miss rate". In n -AFC paradigms γ usually reflects chance performance and is fixed at the reciprocal of the number of intervals per trial, and λ reflects the stimulus-independent error rate or "lapse rate"—see Wichmann and Hill (2000) for more details.

When a parameter set has been estimated, we will usually be interested in measurements of the threshold (displacement along the x -axis) and slope of the psychometric function. We calculate thresholds by taking the inverse of F at a specified probability level, usually 0.5. Slopes are calculated by finding the derivative of F with respect to x , evaluated at a specified threshold. Thus we shall use the notation $\text{threshold}_{0.8}$, for example, to mean $F_{0.8}^{-1}$, and $\text{slope}_{0.8}$ to mean dF/dx evaluated at $F_{0.8}^{-1}$. When we use the terms "threshold" and "slope" without a subscript we mean $\text{threshold}_{0.5}$ and $\text{slope}_{0.5}$: in our 2-AFC examples this will mean the stimulus value and slope of F at the point where performance is approximately 75% correct, although the exact performance level is affected slightly by the (small) value of λ .

Where an estimate of a parameter set is required, given a particular data set, we use a maximum-likelihood search algorithm, with Bayesian constraints on the parameters based on our beliefs about their possible values. For example, λ is constrained within the range $[0, 0.06]$, reflecting our belief that normal, trained observers do not make stimulus-independent errors at high rates. We describe our method in detail in Wichmann and Hill (2000).

Estimates of variability: asymptotic versus Monte Carlo methods

In order to be able to compare response thresholds or slopes across experimental conditions, experimenters require a measure of their variability, which will depend on the number of experimental trials taken and their placement along the stimulus axis. Thus a fitting procedure must not only provide parameter estimates, but also error estimates for those parameters. Reporting error estimates on fitted parameters is unfortunately not very common in psychophysical studies. Sometimes Probit Analysis has been used to provide variability estimates (Finney, 1952; Finney, 1971). In Probit Analysis an iteratively

reweighted linear regression is performed on the data once they have undergone transformation through the inverse of a cumulative Gaussian function. Probit Analysis relies, however, on asymptotic theory: maximum-likelihood estimators are asymptotically Gaussian, allowing the standard deviation to be computed from the empirical distribution (Cox & Hinkley, 1974). Asymptotic methods assume that the number of data points is large—unfortunately, however, the number of points in a typical psychophysical data set is small (between 4 and 10, with between 20 and 100 trials at each) and in these cases, substantial errors have accordingly been found in the probit estimates of variability (Foster & Bischof, 1987; Foster & Bischof, 1991; McKee, Klein, & Teller, 1985). For this reason asymptotic theory methods are not recommended for estimating variability in most realistic psychophysical settings.

An alternative method, the *bootstrap* (Efron, 1979; Efron, 1982; Efron & Gong, 1983; Efron & Tibshirani, 1991; Efron & Tibshirani, 1993), has been made possible by the recent sharp increase in the processing speed of desktop computers. The bootstrap method is a Monte Carlo resampling technique relying on a large number of simulated repetitions of the original experiment. It is potentially well suited to the analysis of psychophysical data because its accuracy does not rely on large numbers of trials as do methods derived from asymptotic theory (Hinkley, 1988). We apply the bootstrap to the problem of estimating the variability of parameters, thresholds and slopes of psychometric functions, following Maloney (1990), Foster and Bischof (1987; 1991; 1997) and Treutwein (1995; Treutwein & Strasburger, 1999).

The essence of Monte Carlo techniques is that a large number, B , of “synthetic” data sets $\mathbf{y}_1^*, \dots, \mathbf{y}_B^*$ are generated. For each data set \mathbf{y}_i^* , the quantity of interest ϑ (threshold or slope, for example) is estimated to give $\hat{\vartheta}_i^*$. The process for obtaining $\hat{\vartheta}_i^*$ is the same as that used to obtain the first estimate $\hat{\vartheta}$. Thus, if our first estimate was obtained by $\hat{\vartheta} = t(\hat{\boldsymbol{\theta}})$, where $\hat{\boldsymbol{\theta}}$ is the maximum-likelihood parameter estimate from a fit to the orig-

inal data \mathbf{y} , so the simulated estimates $\hat{\vartheta}_i^*$ will be given by $\hat{\vartheta}_i^* = t(\hat{\boldsymbol{\theta}}_i^*)$, where $\hat{\boldsymbol{\theta}}_i^*$ is the maximum-likelihood parameter estimate from a fit to the simulated data \mathbf{y}_i^* .

Sometimes it is erroneously assumed that the intention is to measure the variability of the underlying ϑ itself. This cannot be the case, however, because repeated computer simulation of the same experiment is no substitute for the real repeated measurements this would require. What Monte Carlo simulations *can* do is estimate the variability inherent in (i) our sampling as characterized by the distribution of sample points (\mathbf{x}) and the size of the samples (\mathbf{n}), and (ii) any interaction between our sampling strategy and the process used to estimate ϑ : i.e. assuming a model of the observer's variability, fitting a function to obtain $\hat{\boldsymbol{\theta}}$ and applying $t(\hat{\boldsymbol{\theta}})$.

Bootstrap data sets: non-parametric and parametric generation

In applying Monte Carlo techniques to psychophysical data, we require, in order to obtain a simulated data set \mathbf{y}_i^* , some system that provides generating probabilities \mathbf{p} for the binomial variates $y_{i1}^*, \dots, y_{iK}^*$. These should be the same generating probabilities that we hypothesize to underlie the empirical data set \mathbf{y} .

Efron's bootstrap offers such a system. In the non-parametric bootstrap method we would assume $\mathbf{p}=\mathbf{y}$. This is equivalent to resampling, with replacement, the original set of correct and incorrect responses on each block of observations j in \mathbf{y} to produce a simulated sample y_{ij}^* .

Alternatively, a parametric bootstrap can be performed. In the parametric bootstrap assumptions are made about the generating model from which the observed data are believed to arise. In the context of estimating the variability of parameters of psychometric functions, the data are generated by a simulated observer whose underlying probabilities of success are determined by the maximum-likelihood fit to the real observer's data ($\mathbf{y}_{fit} = \psi(\mathbf{x}; \hat{\boldsymbol{\theta}})$). Thus where the non-parametric bootstrap uses \mathbf{y} , the parametric bootstrap uses \mathbf{y}_{fit} as generating probabilities \mathbf{p} for the simulated data sets.

As is frequently the case in statistics, the choice of parametric versus non-parametric analysis concerns how much confidence one has in one's hypothesis about the underlying mechanism that gave rise to the raw data as against the confidence one has in the raw data's precise numerical values. Choosing the parametric bootstrap for the estimation of variability in psychometric function fitting appears the natural choice for several reasons. First and foremost, in fitting a parametric model to the data one has already committed oneself to a parametric analysis. No additional assumptions are required to perform a parametric bootstrap beyond those required for fitting a function to the data: specification of the source of variability (binomial variability) and the model from which the data are most likely to come (parameter vector θ and distribution function F). Second, given the assumption that data from psychophysical experiments are generated from Bernoulli processes, we *expect* data to be variable ("noisy"). The non-parametric bootstrap treats every data point as if its exact value reflected the underlying mechanism¹. The parametric bootstrap, on the other hand, allows the data points to be treated as noisy samples from a smooth and monotonic function, determined by θ and F .

One consequence of the two different bootstrap regimes is as follows: Assume two observers performing the same psychophysical task at the same stimulus intensities \mathbf{x} , and assume that it happens that the maximum-likelihood fits to the two data sets yield identical parameter vectors θ . Given such a scenario, the parametric bootstrap returns identical estimates of variability for both observers, since it only depends on \mathbf{x} , θ and F . The non-parametric bootstrap's estimates would, on the other hand, depend on the individual differences between the two data sets \mathbf{y}_1 and \mathbf{y}_2 —something we consider unconvincing: a method for estimating variability in parameters and thresholds should return identical estimates for identical observers performing the identical experiment.²

¹Under different circumstances and in the absence of a model of the noise and/or the process from which the data stem, this is frequently the best one can do.

Treutwein (1995) and Treutwein and Strasburger (1999) used the non-parametric bootstrap, and Maloney (1990) used the parametric bootstrap, to compare bootstrap estimates of variability with real-word variability in the data of repeated psychophysical experiments. All of the above studies found bootstrap studies to be in agreement with the human data. Keeping in mind that the number of repeats in the above quoted cases was small, this is nonetheless encouraging, suggesting that bootstrap methods are a valid method of variability estimation for parameters fitted to psychophysical data.

Testing the bridging assumption

Asymptotically, that is for large K and N , $\hat{\theta}$ will converge towards θ since maximum-likelihood estimation is asymptotically unbiased³ (Cox & Hinkley, 1974; Kendall & Stuart, 1979). For the small K typical of psychophysical experiments, however, we can only hope that our estimated parameter vector $\hat{\theta}$ is “close enough” to the true parameter vector θ for the estimated variability in the parameter vector $\hat{\theta}$ obtained by the bootstrap method to be valid. We call this the bootstrap bridging assumption.

Whether $\hat{\theta}$ is indeed sufficiently close to θ depends, in a complex way, on the sampling, that is the number of blocks of trials (K), the numbers of observations at each block of trials (\mathbf{n}), and the stimulus intensities (\mathbf{x}) relative to the true parameter vector θ . Maloney (1990) summarized these dependencies for a given experimental design by plotting the standard deviation of $\hat{\beta}$ as a function of α and β as a contour plot (Maloney, 1990, Fig. 3, p. 129). Similar contour plots for the standard deviation of $\hat{\alpha}$, and for bias in both $\hat{\alpha}$ and $\hat{\beta}$, could be obtained. If, due to small K , bad sampling, or otherwise, our

²This is only true, of course, if we have reason to believe that our model actually is a good model of the process under study. This important issue is taken up in the section on goodness-of-fit in our companion paper.

³This only holds if our model is correct: maximum-likelihood parameter estimation for two-parameter psychometric functions to data from observers who occasionally lapse, i.e. display non-stationarity, is asymptotically *biased*, as we show in our companion paper, together with a method to overcome such bias (Wichmann & Hill, 2000).

estimation procedure is inaccurate, the distribution of bootstrap parameter vectors $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$ —centred around $\hat{\theta}$ —will not be centred around true θ . As a result the estimates of variability are likely to be incorrect unless the magnitude of the standard deviation is similar around $\hat{\theta}$ and θ despite the fact that the points are some distance apart in parameter space.

One way to assess the likely accuracy of bootstrap estimates of variability is to follow Maloney and to examine the local flatness of the contours around $\hat{\theta}$, our best estimate of θ . If the contours are sufficiently flat then the variability estimates will be similar, assuming that the true θ is somewhere within this flat region. However, the process of local contour estimation is computationally expensive as it requires a very large number of complete bootstrap runs for each data set.

A much quicker alternative is the following: Having obtained $\hat{\theta}$, and performed a bootstrap using $\psi(x; \hat{\theta})$ as the generating function, we move to eight different points ϕ_1, \dots, ϕ_8 in α - β space. Eight Monte Carlo simulations are performed, using ϕ_1, \dots, ϕ_8 as the generating functions to explore the variability in those parts of the parameter space (only the generating parameters of the bootstrap are changed: \mathbf{x} remains the same for all of them). If the contours of variability around $\hat{\theta}$ are sufficiently flat, as we hope they are, then confidence intervals at ϕ_1, \dots, ϕ_8 should be of the same magnitude as those obtained at $\hat{\theta}$. Prudence should lead us to accept the largest of the nine confidence intervals obtained as our estimate of variability.

A decision has to be made as to which eight points in α - β space to use for the new set of simulations. Generally, provided that the psychometric function is at least reasonably well sampled, the contours vary smoothly in the immediate vicinity of θ , so that the precise placement of the sample points ϕ_1, \dots, ϕ_8 is not critical. One suggested and easy way to obtain a set of additional generating parameters is shown in Fig. 1.

— Insert Figure 1 here —

Figure 1 shows $B=2000$ bootstrap parameter pairs as dark filled circles plotted in α - β space. Simulated data sets were generated from ψ_{gen} with the Weibull as F and $\theta_{gen} = \{10, 3, 0.5, 0.01\}$; the sampling scheme no. 7 (triangles; see Fig. 2) was used and N set to 480 ($n_i=80$). The large central triangle at (10, 3) marks the generating parameter set; the solid and dashed line segments adjacent to the x - and y -axes mark the 68% and 95% confidence intervals for α and β , respectively⁴. In the following we shall use WCI to stand for width of confidence interval, with a subscript denoting its coverage percentage, i.e. WCI_{68} denotes the width of the 68 % confidence interval⁵. The eight additional generating parameter pairs ϕ_1, \dots, ϕ_8 are marked by the light triangles. They form a rectangle whose sides have length WCI_{68} in α and β . Typically, this central rectangular region contains approximately 50% of all α - β pairs and could thus be viewed as a crude joint 50%-confidence region for α and β . A coverage percentage of 50% appears to us a sensible compromise between erroneously accepting the estimate around $\hat{\theta}$, potentially underestimating the true variability, and performing additional bootstrap replications too far in the periphery where variability becomes erroneously inflated due to poor sampling. In terms of statistical hypothesis testing, we try to balance Type I and II errors: small error bars increase Type I errors (falsely rejecting a null-hypothesis, H_0 , that there is no difference between two experimental conditions), inflated error bars decrease the power of the test and increase Type II errors (failing to reject H_0).

Monte Carlo Simulations

In both our papers we use only the specific case of the 2-AFC paradigm in our examples: thus γ is fixed at 0.5. In our simulations, where we must assume a distribution of “true”

⁴Confidence intervals here are computed by the bootstrap percentile method: the 95% confidence interval for α , for example, was determined simply by $[\alpha^{*(0.025)}, \alpha^{*(0.975)}]$, where $\alpha^{*(n)}$ denotes the 100n-th percentile of the bootstrap distribution α^* .

⁵68% was chosen because this is the approximate coverage of the familiar “standard error bar” denoting one’s original estimate \pm one standard deviation of a Gaussian.

generating probabilities, we always use the Weibull function in conjunction with the same fixed set of generating parameters $\theta_{\text{gen}}: \{\alpha_{\text{gen}} = 10, \beta_{\text{gen}} = 3, \gamma_{\text{gen}} = 0.5, \lambda_{\text{gen}} = 0.01\}$. In our investigation of the effects of sampling patterns we shall always use $K=6$ and n_i constant, that is, 6 blocks of trials with the same number of points in each block. The number of observations per point, n_p , could be set to 20, 40, 80 or 160, and with $K=6$, this means that the total number of observations N could take the values 120, 240, 480 and 960.

We have introduced these limitations purely for the purposes of illustration, to keep our explanatory variables down to a manageable number. We have found, in many other simulations, that in in most cases this is done without loss of generality of our conclusions.

The effects of sampling schemes and number of trials

One of our aims in this study was to examine the effect of N and one's choice of sample points \mathbf{x} on both the size of one's confidence intervals for $\hat{\theta}$ and their sensitivity to errors in $\hat{\theta}$.

Seven different sampling schemes were used, each dictating a different distribution of data points along the stimulus axis; they are the same schemes as used and described in Wichmann and Hill (2000), and they are shown in Fig. 2. Each horizontal chain of symbols represents one of the schemes, marking the stimulus values at which the six sample points are placed. The different symbol shapes will be used to identify the sampling schemes in our results plots. To provide a frame of reference, the solid curve shows the psychometric function used, i.e. $0.5 + 0.5F(x; \{\alpha_{\text{gen}}, \beta_{\text{gen}}\})$, with the 55%, 75% and 95% performance levels marked by dotted lines.

— Insert Figure 2 here —

As we shall see, even for a fixed number of sample points and a fixed number of trials per point, biases in parameter estimation and goodness-of-fit assessment (companion

paper) as well as the width of confidence intervals (this paper), all depend markedly on the distribution of stimulus values \mathbf{x} .

Monte Carlo data sets were generated using our seven sampling schemes shown in Fig. 2 using the generation parameters $\boldsymbol{\theta}_{gen}$, as well as N and K as specified above. A maximum-likelihood fit was performed on each simulated data set to obtain bootstrap parameter vectors $\hat{\boldsymbol{\theta}}^*$ from which we subsequently derived the x -values corresponding to $threshold_{0.5}$ and $threshold_{0.8}$ as well as to the slope. For each sampling scheme and value of N a total of nine simulations were performed: one at $\boldsymbol{\theta}_{gen}$ and eight more at points, $\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_8$, as specified in our section on the bootstrap bridging assumption. Thus, each of our 28 conditions (7 sampling schemes \times 4 values of N) required 9×2000 simulated datasets, for a total of 304,000 simulations, or 2.268×10^8 simulated 2-AFC trials.

Figures 3, 4 and 5 show the results of the simulations dealing with $slope_{0.5}$, $threshold_{0.5}$ and $threshold_{0.8}$, respectively. The left panel of each figure plots the WCI_{68} of the estimate under consideration as a function of N . Data for all seven sampling schemes are shown using their respective symbols. The right hand panel plots, as a function of N , the maximal elevation of the WCI_{68} encountered in the vicinity of $\boldsymbol{\theta}_{gen}$, that is $\max\{WCI_{\boldsymbol{\phi}_1}/WCI_{\boldsymbol{\theta}_{gen}}, \dots, WCI_{\boldsymbol{\phi}_8}/WCI_{\boldsymbol{\theta}_{gen}}\}$. The elevation factor is an indication of the sensitivity of our variability estimates to errors in the estimation of $\boldsymbol{\theta}$. The smaller it is, the better.

— Insert Figure 3 here —

The left panel of Fig. 3 shows the WCI_{68} around the median estimated slope. Clearly, the different sampling schemes have a profound effect on the magnitude of the confidence intervals for slope estimates. For example, in order to ensure that the WCI_{68} is approximately 0.06, one requires nearly 960 trials if using sampling s1 or s4. Sampling schemes s3 or s7, on the other hand, require only around 300 trials to achieve similar

confidence interval width. The important difference that makes s_3 and s_7 more efficient than s_1 and s_4 is the presence of samples at high predicted performance values ($p \geq 0.9$) where binomial variability is low and thus the data constrain our maximum-likelihood fit more tightly. The right panel of Fig. 3 illustrates the complex interactions between different sampling schemes, N , and the stability of the bootstrap estimates of variability as indicated by the local flatness of the contours around θ_{gen} . A perfectly flat local contour would result in the horizontal line at 1. Sampling scheme s_3 is well-behaved for $N \geq 240$, its maximal elevation being below 1.5. For $N = 120$, however, elevation rises sharply to near 3. Similarly, s_5 is well-behaved for $N \geq 480$, but its bootstrap estimate of the WCI_{68} around the slope for $N = 120$ is extremely unreliable. Other schemes, like s_1 , s_2 , or s_6 , never rise above an elevation of 2 regardless of N . It is important to note that the magnitude of WCI_{68} at θ_{gen} (left panel) is by no means a good predictor of the *stability* of the estimate, as indicated by the sensitivity factor (right panel).

— Insert Figure 4 here —

— Insert Figure 5 here —

Figure 4, showing the equivalent data for the estimates of $threshold_{0.5}$, also illustrates that some sampling schemes make much more efficient use of the experimenter's time by providing WCI_{68} 's a factor of 1.6 more compact than others.

Two aspects of the data shown in the right panel of Fig. 4 are important. First, the sampling schemes fall into two distinct classes: five of the seven sampling schemes are almost ideally well-behaved with elevations barely reaching 1.4 even for small N ; the remaining two, on the other hand, behave poorly with elevations in excess of 3 for $N = 240$ or below. The two unstable sampling schemes, s_1 and s_4 , are those that do *not* include at least a single sample point at $p \geq 0.95$ (see Fig. 2). It thus appears crucial to include at least one sample point at $p \geq 0.95$ to make one's estimates of variability robust to small errors in $\hat{\theta}$, even if the threshold of interest, as in our example, has a p of only approximately 0.75.

Second, s_1 is prone to lead to Type I errors if the sensitivity analysis (or bootstrap bridging assumption test) is not carried out: the WCI_{68} is unstable in the vicinity of θ_{gen} even though WCI_{68} is small at θ_{gen} .

Figure 5 is similar to Fig. 4 except that it shows the WCI_{68} around $threshold_{0.8}$. The trends found in Fig. 4 are even more exaggerated here: The sampling schemes without high sampling points (s_1, s_4) are unstable to the point of being meaningless for $N < 480$. In addition, their WCI_{68} is inflated relative to that of the other sampling schemes even at θ_{gen} .

The results of the Monte Carlo simulations are summarized in Table 1. The columns of Table 1 correspond to the different sampling schemes, marked by their respective symbols. The first four rows contain the WCI_{68} at $threshold_{0.5}$ for $N = 120, 240, 480$ and 960 ; similarly, the next four rows contain the WCI_{68} at $threshold_{0.8}$ and the following four those for the $slope_{0.5}$. Each entry in the table corresponds to the largest of the WCI_{68} 's measurements at θ_{gen} and ϕ_1, \dots, ϕ_8 for a given sampling scheme and N ; in terms of Figures 3, 4 and 5, it is the product of the value for a given WCI_{68} in the left panel multiplied by the appropriate elevation factor of its right panel. This quantity we denote as $MWCI_{68}$, standing for maximum width of the 68% confidence interval. The scheme with the lowest $MWCI_{68}$ in each row is given the score 100. The others on the same row are given proportionally higher scores, to indicate their $MWCI_{68}$ as a percentage of the best scheme's value. The bottom three rows of Table 1 contain summary statistics of how well the different sampling schemes do across all twelve estimates.

— Insert Table 1 here —

Inspecting Table 1 reveals that the sampling schemes fall into three categories. By a long way worst are sampling schemes s_1 and s_4 with mean and median $MWCI_{68} > 200\%$. Second come three sampling schemes with medians and means between 120% and 140%— s_2, s_5 and s_6 . Each of these has at least one sample at $p \geq 0.95$. s_6 clearly demonstrates the importance of this high sample point. Comparing s_1 and s_6 , s_6 is identical to scheme s_1

except that one sample point was moved from 0.75 to 0.95. Still, s_6 is superior to s_1 on each of the twelve estimates, and often very markedly so. Finally, there are two sampling schemes with medians below 106% and means around 110%—very nearly optimal⁶ on most estimates. Both these sampling schemes, s_3 and s_7 , have 50% of their sample points at $p \geq 0.90$ and one third at $p \geq 0.95$.

In order to obtain stable estimates of variability of parameters, thresholds and slopes of psychometric functions, it appears that we must include at least one, but preferably more, sample points at large p values. Such sampling schemes are, however, sensitive to stimulus independent lapses that could potentially bias the estimates if we were to fix the upper asymptote of the psychometric function (the parameter λ in equation (1); see our companion paper, Wichmann & Hill, 2000).

Somewhat counter-intuitively, it is thus not sensible to place all or most samples close to the point of interest (for example close to $\text{threshold}_{0.5}$, in order to obtain tight confidence intervals for $\text{threshold}_{0.5}$), because estimation is done via the the whole psychometric function which in turn is estimated from the entire dataset. Hence adaptive techniques that sample predominantly around the threshold value of interest, are less efficient than one might think (c.f. Lam, Mills, & Dubno, 1996).

Influence of the distribution function on estimates of variability

Thus far we have argued in favour of the bootstrap method for estimating the variability of fitted parameters, thresholds and slopes, since its estimates do not rely on asymptotic theory. However, in the context of fitting psychometric functions one requires in addition that the exact form of the distribution function F —Weibull, logistic, cumulative Gaussian, Gumbel or any other reasonably similar sigmoid—has only a minor influence on the estimates of variability. The importance of this cannot be underestimated since a strong dependence of the estimates of variability on the precise algebraic form of the dis-

⁶Optimal here of course means relative to the sampling schemes explored.

tribution function would call the usefulness of the bootstrap into question because, as experimenters, we do not know, and never will, the true underlying distribution function or objective function from which the empirical data were generated. The problem is illustrated in Fig. 6; Figure 6(a) shows four different psychometric functions: 1.) $\psi_W(x; \boldsymbol{\theta}_W)$, using the Weibull as F , and $\boldsymbol{\theta}_W = \{10, 3, 0.5, 0.01\}$ (our “standard” generating function ψ_{gen}). 2.) $\psi_{CG}(x; \boldsymbol{\theta}_{CG})$, using the cumulative Gaussian with $\boldsymbol{\theta}_{CG} = \{8.875, 3.278, 0.5, 0.01\}$. 3.) $\psi_L(x; \boldsymbol{\theta}_L)$, using the logistic with $\boldsymbol{\theta}_L = \{8.957, 2.014, 0.5, 0.01\}$ and, finally, 4.) $\psi_G(x; \boldsymbol{\theta}_G)$, using the Gumbel and $\boldsymbol{\theta}_G = \{10.022, 2.906, 0.5, 0.01\}$. For all practical purposes in psychophysics the four functions are indistinguishable. Thus, if one of the above psychometric functions were to provide a good fit to a data set, all of them would despite that at most one of them is correct. The question one has to ask is whether making the choice of one distribution function over another markedly changes the bootstrap estimates of variability.⁷ Note that this is not trivially true: whilst it can be the case that several psychometric functions with different distribution functions F are indistinguishable *given* a particular dataset—as shown in Fig. 6(a)—this does not imply that the same is true for every dataset *generated from* one of such similar psychometric functions during the bootstrap procedure: Figure 6(b) shows the fit of two psychometric functions (Weibull and logistic) to a dataset generated from our “standard” generating function ψ_{gen} , using sampling scheme s2 with $N = 120$.

— Insert Figure 6 here —

Slope, $\text{threshold}_{0.8}$ and $\text{threshold}_{0.2}$ are quite dissimilar for the two fits, illustrating the point that there is a real possibility that the bootstrap distributions of thresholds and

⁷Of course there might be situations where one psychometric function using a particular distribution function provides a significantly better fit to a given data set than others using different distribution functions. Differences in bootstrap estimates of variability in such cases are not worrisome: The appropriate estimates of variability are those of the best fitting function.

slopes from the B bootstrap repeats differ substantially for different choices of F , even if the fits to the original (empirical) dataset were almost identical.

To explore the effect of F on estimates of variability we conducted Monte Carlo simulations using $\Psi = \frac{1}{4}[\Psi_W + \Psi_{CG} + \Psi_L + \Psi_G]$ as generating function, and fitted psychometric functions using the Weibull, cumulative Gaussian, logistic and Gumbel as distribution function to each dataset. From the fitted psychometric functions we obtained estimates of $\text{threshold}_{0.2}$, $\text{threshold}_{0.5}$, $\text{threshold}_{0.8}$ and $\text{slope}_{0.5}$ as described previously. All four different values of N and our seven sampling schemes were used, resulting in 112 conditions (4 distribution functions x 7 sampling schemes x 4 N values). In addition, we repeated the above procedure 40 times to obtain an estimate of the numerical variability intrinsic to our bootstrap routines⁸, for a total of 4,480 bootstrap repeats affording 8,960,000 psychometric function fits (4.032×10^9 simulated 2AFC trials).

An analysis of variance (ANOVA) was applied to the resulting data, with the number of trials N , the sampling schemes s_1 to s_7 , and the distribution function F as independent factors (variables). The dependent variables were the confidence interval widths (WCI_{68}); each cell contained the WCI_{68} estimates from our 40 repetitions. For all four dependent measures— $\text{threshold}_{0.2}$, $\text{threshold}_{0.5}$, $\text{threshold}_{0.8}$ and $\text{slope}_{0.5}$ —not only the first two factors, the number of trials N and the sampling scheme, were, as expected, significant ($p < 0.0001$), but also the distribution function F and all possible interactions: the three two-way interactions and the three-way interaction were similarly significant at $p < 0.0001$. This result in itself, however, is not necessarily damaging to the bootstrap method applied to psychophysical data because the significance is brought about by the very low (and desirable) variability of our WCI_{68} estimates: model R^2 is between 0.995 and 0.997, implying that virtually all the variance in our simulations is due to N , sampling scheme, F and interactions thereof.

⁸ WCI_{68} estimates were obtained using bias corrected and accelerated (BCa) confidence intervals, described in the next section.

Rather than exclusively focusing on significance, in table 2 we provide information about effect size, namely the percentage of the total sum of squares of variation accounted for by the different factors and their interactions. For $\text{threshold}_{0.5}$ and $\text{slope}_{0.5}$ (columns 2 and 4) N , sampling scheme, and their interaction account for 98.63 and 96.39 % of the total variance, respectively⁹. The choice of distribution function F does not have, despite being a significant factor, a large effect on WCI_{68} for $\text{threshold}_{0.5}$ and $\text{slope}_{0.5}$.

— Insert Table 2 here —

The same is not true, however, for the WCI_{68} of $\text{threshold}_{0.2}$. Here the choice of F has an undesirably large effect on the bootstrap estimate of WCI_{68} —its influence is larger than that of the sampling scheme used—and only 84.36 % of the variance is explained by N , sampling scheme, and their interaction. Figure 7, finally, summarizes the effect sizes of N , sampling scheme and F graphically: Each of the four panels of Fig. 7 plots the WCI_{68} (normalized by dividing each WCI_{68} score by the largest mean WCI_{68}) on the y-axis as a function of N on the x-axis; the different symbols refer to the different sampling schemes. The two symbols shown in each panel correspond to the sampling schemes which yielded the smallest and largest mean WCI_{68} (averaged across F and N). The gray levels, finally, code the smallest (black), mean (gray) and largest (white) WCI_{68} for a given N and sampling scheme as a function of the distribution function F .

— Insert Figure 7 here —

For $\text{threshold}_{0.5}$ and $\text{slope}_{0.5}$ (Fig. 7b and d) estimates are virtually unaffected by the choice of F , but for $\text{threshold}_{0.2}$ the choice of F has a profound influence on WCI_{68} (for

⁹Clearly, the above reported effect sizes are tied to the ranges in the factors explored: N spanned a comparatively large range of 120 to 960 observations, or a factor of 8, whereas all of our sampling schemes were “reasonable”—inclusion of “unrealistic” or “unusual” sampling schemes, e.g. all x values such that nominal y values are below 0.55, would have increased the percentage of variation accounted for by sampling scheme. Taken together, N and sampling scheme should be representative of most typically used psychophysical settings, however.

example, in Fig. 7a, there is a difference of nearly a factor of two for sampling scheme s7 (triangles) when $N = 120$). The same is also true, albeit to a lesser extent, if one is interested in $\text{threshold}_{0.8}$: Figure 7c shows the (again undesirable) interaction between sampling scheme and choice of F . WCI_{68} estimates for sampling scheme s5 (leftward triangles) show little influence of F , but for sampling scheme s4 (rightward triangles) the choice of F has a marked influence on WCI_{68} . It was generally the case for $\text{threshold}_{0.8}$ that those sampling schemes that resulted in small confidence intervals (s2, s3, s5, s6 and s7, see previous section) were less affected by F than those resulting in large confidence intervals (s1 and s4).

Two main conclusions can be drawn from these simulations: First, in the absence of any other constraints experimenters should choose as “threshold” and “slope” measures corresponding to $\text{threshold}_{0.5}$ and $\text{slope}_{0.5}$, because only then the main factors influencing the estimates of variability are the number and placement of stimuli, as we would like it to be. Second, away from the midpoint of F estimates of variability are, however, not as independent of the distribution function chosen as one might wish, in particular for lower proportions correct ($\text{threshold}_{0.2}$ is much more affected by the choice of F than $\text{threshold}_{0.8}$, c.f. Fig.’s 7a and c). If very low (or, perhaps, very high) response thresholds must be used when comparing experimental conditions, e.g. 60 % (or 90 %) correct in 2AFC, and only small differences exist between the different experimental conditions, this requires the exploration of a number of distribution functions F to avoid finding significant differences between conditions due to the (arbitrary) choice of a distribution function resulting in comparatively narrow confidence intervals.

Bootstrap confidence intervals

In the existing literature on bootstrap estimates of the parameters and thresholds of psychometric functions, most studies use parametric or non-parametric *plug-in estimates*¹⁰ of the variability of a distribution $\hat{\mathfrak{D}}^*$. For example, Foster & Bischof (1997) estimate parametric (moment-based) standard deviations σ by the plug-in estimate $\hat{\sigma}$. Maloney

(1990), in addition to $\hat{\sigma}$, uses a comparable non-parametric estimate, obtained by scaling plug-in estimates of the interquartile range (IQR) so as to cover a confidence interval of 68.3%. Neither kind of plug-in estimate is guaranteed to be reliable, however: moment-based estimates of a distribution's central tendency (such as the mean) or variability (such as $\hat{\sigma}$) are not robust; they are very sensitive to outliers because a change in a single sample can have an arbitrarily large effect on the estimate (the estimator is said to have a breakdown of $1/n$, because that is the proportion of the data set that can have such an effect. Non-parametric estimates are usually much less sensitive to outliers and the median, for example, has a breakdown of $1/2$ as opposed to $1/n$ for the mean). A moment-based estimate of a quantity ϑ might be seriously in error if only a single bootstrap estimate $\hat{\vartheta}_i^*$ is wrong by a large amount. Large errors can and do occur occasionally, for example, when the maximum-likelihood search algorithm gets stuck in a local minimum on its error surface¹¹.

Non-parametric plug-in estimates are also not without problems. Percentile-based bootstrap confidence intervals are sometimes significantly biased and converge slowly to the true confidence intervals (Efron & Tibshirani, 1993, ch. 12-14, 22). In the psychological literature this problem was critically noted by Rasmussen (1987; 1988).

Methods to improve convergence accuracy and avoid bias have received the most theoretical attention in the study of the bootstrap (Efron, 1987; 1988; Efron & Tibshirani, 1993; Hall, 1988; Hinkley, 1988; Strube, 1988; c.f. Foster & Bischof, 1991, p. 158).

¹⁰A straightforward way to estimate a quantity ϑ which is derived from a probability distribution F by $\vartheta=t(F)$, is to obtain \hat{F} from empirical data and then use $\hat{\vartheta}=t(\hat{F})$ as an estimate. This is called a plug-in estimate.

¹¹Foster & Bischof (1987) report problems with local minima, which they overcame by discarding bootstrap estimates that were larger than 20 times the stimulus range (4.2% of their data points had to be removed). Non-parametric estimates naturally avoid having to perform post-hoc data smoothing by their resilience to such infrequent but extreme outliers.

In situations where asymptotic confidence intervals are known to apply and are correct, *bias-corrected and accelerated* (BC_a) confidence intervals have been demonstrated to show faster convergence and increased accuracy over ordinary percentile based methods while retaining the desirable property of robustness. (See in particular Efron & Tibshirani, 1993, p. 183, Table 14.2 and p. 184, Fig. 14.3, as well as Efron (1988), Rasmussen (1987; 1988) and Strube (1988)).

BC_a confidence intervals are necessary because the distribution of sampling points \mathbf{x} along the stimulus axis may cause the bootstrap estimates $\hat{\boldsymbol{\theta}}^*$ to be biased and skewed estimators of the generating values $\hat{\boldsymbol{\theta}}$. The same applies to the bootstrap distributions of estimates $\hat{\boldsymbol{\theta}}^*$ of any quantity of interest, be it thresholds, slopes or whatever. Maloney found skew and bias particularly problematic for the distribution of the β parameter of the Weibull, $\boldsymbol{\beta}^*$ ($N=210, K=7$). We also found in our simulations that $\boldsymbol{\beta}^*$ —and thus slopes \mathbf{s}^* —were skewed and biased for N smaller than 480, even using the best of our sampling schemes. The BC_a method attempts to correct both bias and skew by assuming that an increasing transformation, m , exists to transform the bootstrap distribution into a normal distribution. Hence we assume $\Phi = m(\theta)$ and $\hat{\Phi} = m(\hat{\theta})$ resulting in

$$\frac{\hat{\Phi} - \Phi}{k_{\Phi}} \sim N(-z_0, 1), \quad (2)$$

where

$$k_{\Phi} = k_{\Phi_0} + a(\Phi - \Phi_0), \quad (3)$$

and k_{Φ_0} any reference point on the scale of Φ values. In equation (3) z_0 is the *bias correction* term, and a in equation (4) is the *acceleration* term. Assuming equation (2) to be correct, it has been shown that an ϵ -level confidence interval endpoint of the BC_a interval can be calculated as

$$\hat{\theta}_{BC_a}[\epsilon] = \hat{G}^{-1} \left(CG \left(\hat{z}_0 + \frac{\hat{z}_0 + z^{(\epsilon)}}{1 - \hat{a}(\hat{z}_0 + z^{(\epsilon)})} \right) \right), \quad (4)$$

with CG being the cumulative Gaussian distribution function, \hat{G}^{-1} is the inverse of the cumulative distribution function of the bootstrap replications $\hat{\theta}^*$, \hat{z}_0 is our estimate of the bias, and \hat{a} is our estimate of acceleration. For details on how to calculate the bias and the acceleration terms see Efron & Tibshirani (1993), ch. 22 and also Davison & Hinkley (1997), ch. 5.

For large classes of problems it has been shown that equation (2) is approximately correct, and that the error in the confidence intervals obtained from equation (4) are smaller than those introduced by the standard percentile approximation to the true underlying distribution, where an ϵ -level confidence interval endpoint is simply $\hat{G}[\epsilon]$ (Efron & Tibshirani, 1993). While we cannot offer a formal proof that this is also true for the bias and skew sometimes found in bootstrap estimates from fits to psychophysical data, to our knowledge it has only been shown that BC_a confidence intervals are either superior or equally good in performance to standard percentiles, but not that they perform significantly worse.

Conclusions

In this paper we have given an account of the procedures we use to estimate the variability of fitted parameters and the derived measures such as thresholds and slopes of psychometric functions.

First, we recommend the use of Efron's parametric bootstrap technique, because traditional asymptotic methods have been found to be unsuitable given the small number of data-points typically taken in psychophysical experiments. Second, we have introduced a practicable test of the *bootstrap bridging assumption* or *sensitivity analysis* which must be applied every time bootstrap-derived variability estimates are obtained to ensure that

variability estimates do not change markedly with small variations in the bootstrap generating function's parameters. This is critical because the fitted parameters $\hat{\theta}$ are almost certain to deviate at least slightly from the (unknown) underlying parameters θ . Third, we explored the influence of different sampling schemes (\mathbf{x}) on both the size of one's confidence intervals as well as their sensitivity to errors in $\hat{\theta}$. We conclude that only sampling schemes including at least one sample at $p \geq 0.95$ yield reliable bootstrap confidence intervals. Fourth, we have shown that the size of bootstrap confidence intervals is mainly influenced by \mathbf{x} and N if and only if we choose as threshold and slope values around the midpoint of the distribution function F —particularly for low thresholds ($\text{threshold}_{0.2}$) the precise mathematical form of F exerts a noticeable and undesirable influence on the size of bootstrap confidence intervals. Finally, we have reported that the use of *bias-corrected and accelerated* (BC_a) confidence intervals which improve on parametric and percentile-based bootstrap confidence intervals, whose bias and slow convergence had previously been noted (Rasmussen, 1987).

Together with our companion paper (Wichmann & Hill, 2000) we cover the three central aspects of modelling experimental data: First, parameter estimation, second, obtaining error estimates on these parameters and, third, assessing goodness-of-fit between model and data.

Acknowledgements

For part of this work Felix A. Wichmann was supported by a Wellcome Trust Mathematical Biology studentship and held a Jubilee Scholarship from St. Hugh's College, Oxford. He is now a Fellow by Examination at Magdalen College, Oxford. N. Jeremy Hill is supported by a grant from the Christopher Welch Trust Fund, and a Maplethorpe Scholarship from St. Hugh's College, Oxford.

We are indebted to Andrew Derrington, Karl Gegenfurtner, Bruce Henning, Larry Maloney, Eero Simoncelli and Stefaan Tibeau for helpful comments and suggestions. This paper benefitted considerably from conscientious peer review and we wish to thank our reviewers David Foster, Marjorie Leek and Bernhard Treutwein, as well as the editor, Neil Macmillan, for helping us to improve our manuscript. Part of this work was presented at the Computers in Psychology (CiP) Conference in York, UK, during April 1998 (Hill & Wichmann, 1998). Software implementing the methods described in this paper is available (MATLAB); contact FAW at the address provided or see <http://users.ox.ac.uk/~sruoxfor/psychofit/>

References

- Cox, D. R., & Hinkley, D. V. (1974). Theoretical Statistics. London: Chapman and Hall.
- Davison, A. C., & Hinkley, D. V. (1997). Bootstrap Methods and their Application. Cambridge: CUP.
- Efron, B. (1979). Bootstrap methods: another look at the jackknife. Annals of Statistics, 7, 1-26.
- Efron, B. (1982). The jackknife, the bootstrap and other resampling plans, CBMS-NSF Regional Conference Series in Applied Mathematics . Philadelphia: Society for Industrial and Applied Mathematics.
- Efron, B. (1987). Better bootstrap confidence intervals. Journal of the American Statistical Association, 82, 171-200.
- Efron, B. (1988). Bootstrap confidence intervals: good or bad ? Psychological Bulletin, 104(2), 293-296.
- Efron, B., & Gong, G. (1983). A leisurely look at the bootstrap, the jackknife, and cross-validation. The American Statistician, 37, 36-48.
- Efron, B., & Tibshirani, R. (1991). Statistical data analysis in the computer age. Science, 253, 390-395.
- Efron, B., & Tibshirani, R. J. (1993). An Introduction to the Bootstrap. New York: Chapman & Hall.
- Finney, D. J. (1952). Probit Analysis. (2nd ed.). Cambridge: Cambridge University Press.
- Finney, D. J. (1971). Probit Analysis. (3rd ed.). Cambridge: Cambridge University Press.
- Foster, D. H., & Bischof, W. F. (1987). Bootstrap variance estimators for the parameters of small-sample sensory-performance functions. Biological Cybernetics, 57, 341-347.

- Foster, D. H., & Bischof, W. F. (1991). Thresholds from psychometric functions: superiority of bootstrap to incremental and probit variance estimators. Psychological Bulletin, 109, 152-159.
- Foster, D. H., & Bischof, W. F. (1997). Bootstrap estimates of the statistical accuracy of thresholds obtained from psychometric functions. Spatial Vision, 11(1), 135-139.
- Hall, P. (1988). Theoretical comparison of bootstrap confidence intervals. (with discussion). Annals of Statistics, 16, 927-953.
- Hill, N. J., & Wichmann, F. A. (1998). A bootstrap method for testing hypotheses concerning psychometric functions. Paper presented at the Computers in Psychology, York, UK.
- Hinkley, D. V. (1988). Bootstrap methods. Journal of the Royal Statistical Society B, 50, 321-337.
- Kendall, M. K., & Stuart, A. (1979). The advanced theory of statistics: Vol. 2. Inference and Relationship. New York: Macmillan.
- Lam, C. F., Mills, J. H., & Dubno, J. R. (1996). Placement of observations for the efficient estimation of a psychometric function. Journal of the Acoustical Society of America, 99(6), 3689-3693.
- Maloney, L. T. (1990). Confidence intervals for the parameters of psychometric functions. Perception and Psychophysics, 47, 127-134.
- McKee, S. P., Klein, S. A., & Teller, D. Y. (1985). Statistical properties of forced-choice psychometric functions: Implications of probit analysis. Perception and Psychophysics, 37, 286-298.
- Rasmussen, J. L. (1987). Estimating correlation coefficients: bootstrap and parametric approaches. Psychological Bulletin, 101(1), 136-139.

- Rasmussen, J. L. (1988). Bootstrap confidence intervals: good or bad: Comments on Efron (1988) and Strube (1988) and further evaluation. Psychological Bulletin, 104(2), 297-299.
- Strube, M. J. (1988). Bootstrap type I error rates for the correlation coefficient: an examination of alternate procedures. Psychological Bulletin, 104(2), 290-292.
- Treutwein, B. (1995). Error estimates for the parameters of psychometric functions from a single session. Poster presented at the European Conference of Visual Perception, Tübingen, FRG.
- Treutwein, B., & Strasburger, H. (1999). Assessing the Variability of Psychometric Functions. Paper presented at the European Mathematical Psychology Meeting, Mannheim, FRG.
- Wichmann, F. A., & Hill, N. J. (2000). The psychometric function I: Fitting, sampling and goodness-of-fit. Perception and Psychophysics, (submitted).

Table Captions

Table 1. Columns correspond to the seven sampling schemes and are marked by their respective symbols (see Fig. 2). The first four rows contain the $MWCI_{68}$ at $\text{threshold}_{0.5}$ for $N=120, 240, 480$ and 960 ; similarly, the next eight rows contain the $MWCI_{68}$ at $\text{threshold}_{0.8}$ and $\text{slope}_{0.5}$. (See text for the definition of the $MWCI_{68}$.) Each entry corresponds to the largest $MWCI_{68}$ in the vicinity of θ_{gen} , as sampled at the points θ_{gen} and ϕ_1, \dots, ϕ_8 . The $MWCI_{68}$ values are expressed in percent relative to the minimal $MWCI_{68}$ per row. The bottom three rows of Table 1 contain summary statistics of how well the different sampling schemes perform across estimates.

Table 2. Summary of ANOVA effect size (sum of squares (SS) normalized to 100 percent). The columns refer to $\text{threshold}_{0.2}$, $\text{threshold}_{0.5}$, $\text{threshold}_{0.8}$, and $\text{slope}_{0.5}$, respectively, i.e. to approximately 60, 75 and 90% correct and the slope at 75% correct during 2AFC. Rows correspond to the independent variables, their interactions and summary statistics. See text for details.

Tables

Table 1:








								
		s1	s2	s3	s4	s5	s6	s7
MWCI ₆₈	<i>N</i> = 120	401	127	134	844	136	100	120
at	<i>N</i> = 240	236	121	131	361	135	100	116
$x = F_{0.5}^{-1}$	<i>N</i> = 480	129	103	117	175	123	100	111
	<i>N</i> = 960	109	108	125	140	131	100	115
MWCI ₆₈	<i>N</i> = 120	4137	185	101	2425	112	120	100
at	<i>N</i> = 240	2267	138	101	979	117	131	100
$x = F_{0.8}^{-1}$	<i>N</i> = 480	457	120	100	464	121	138	105
	<i>N</i> = 960	284	105	100	354	120	142	103
MWCI ₆₈	<i>N</i> = 120	147	107	123	260	237	104	100
of dF/dx	<i>N</i> = 240	161	118	100	318	166	135	123
at	<i>N</i> = 480	174	119	100	269	116	153	108
$x = F_{0.5}^{-1}$	<i>N</i> = 960	170	110	100	231	114	167	102
MEAN		723	122	111	568	136	124	109
standard deviation (SD)		1229	22.3	13.8	639	35	24	8
MEDIAN		205	119	101	336	122	126	106

Table 2:

percentage of ANOVA Sum of Squares (SS) accounted for by ...	threshold _{0,2}	threshold _{0,5}	threshold _{0,8}	slope _{0,5}
the number of trials <i>N</i>	76.24	87.30	65.14	72.86
sampling schemes s1 ... s7	6.60	10.00	19.93	19.69
distribution function <i>F</i>	11.36	0.13	3.87	0.92
error (numerical variability in bootstrap)	0.34	0.35	0.46	0.34
interaction of <i>N</i> and sampling schemes s1 ... s7	1.18	0.98	5.97	3.50
sum of interactions involving <i>F</i>	4.28	1.24	4.63	2.69
percentage of SS accounted for without <i>F</i>	84.36	98.63	91.5	96.39

Figure Captions

Figure 1. $B=2000$ datasets were generated from a 2-AFC Weibull psychometric function with parameter vector $\boldsymbol{\theta}_{\text{gen}} = \{10, 3, 0.5, 0.01\}$ and then fit using our maximum-likelihood procedure resulting in 2000 estimated parameter pairs $(\hat{\alpha}, \hat{\beta})$ shown as dark circles in α - β parameter space. The location of the generating α and β (10, 3) is marked by the large triangle in the centre of the plot. The sampling scheme $s7$ was used to generate the datasets (see Fig. 2 for details) with $N=480$. Solid lines mark the 68%-confidence interval width (WCI_{68}) separately for α and β ; broken lines mark the 95%-confidence intervals. The light small triangles show the α - β parameter sets ϕ_1, \dots, ϕ_8 from which each bootstrap is repeated during sensitivity analysis while keeping the x -values of the sampling scheme unchanged.

Figure 2. Shows a 2-AFC Weibull psychometric function with parameter vector $\boldsymbol{\theta} = \{10, 3, 0.5, 0\}$ on semi-logarithmic coordinates. The rows of symbols below the curve mark the x -values of the seven different sampling schemes, $s1$ to $s7$, used throughout the remainder of the paper.

Figure 3. The left hand panel shows the width of the 68%-confidence (WCI_{68}) interval around the median estimate of the distribution of slopes of $B=2000$ fitted psychometric functions to parametric bootstrap datasets generated from $\boldsymbol{\theta}_{\text{gen}} = \{10, 3, 0.5, 0.01\}$ as function of the total number of observations, N . The right hand panel shows the maximal elevation of the WCI_{68} in the vicinity of $\boldsymbol{\theta}_{\text{gen}}$ again as function of N (see text for details). The seven symbols denote the seven sampling schemes as of Fig. 2.

Figure 4. Similar to Fig. 3 except that it shows thresholds corresponding to $F_{0.5}^{-1}$ (approximately equal to 75% correct during 2-AFC).

Figure 5. Similar to Fig. 3 except that it shows thresholds corresponding to $F_{0.8}^{-1}$ (approximately equal to 90% correct during 2-AFC)

Figure 6. (a) Shows four 2–AFC psychometric functions plotted on semi–logarithmic coordinates; each has a different distribution function F (Weibull, cumulative Gaussian, logistic and Gumbel). See text for details. (b) Shows a fit of two psychometric functions with different distribution functions F (Weibull, logistic) to the same dataset, generated from the mean of the four psychometric functions shown in (a) using sampling scheme s_2 with $N = 120$. See text for details.

Figure 7. The four panels show the width of WCI_{68} as a function of N , sampling scheme, as well as the distribution function F . The symbols refer to the different sampling schemes as of Figure 2; gray–levels code the influence of the distribution function F on the width of WCI_{68} for a given sampling scheme and N : black symbols show the smallest WCI_{68} obtained, middle gray the mean WCI_{68} across the four distribution functions used, and white the largest WCI_{68} . (a) shows WCI_{68} for $\text{threshold}_{0.2}$ (b) WCI_{68} for $\text{threshold}_{0.5}$ (c) WCI_{68} for $\text{threshold}_{0.8}$ (d) WCI_{68} for $\text{slope}_{0.5}$.

Figures

Figure 1:

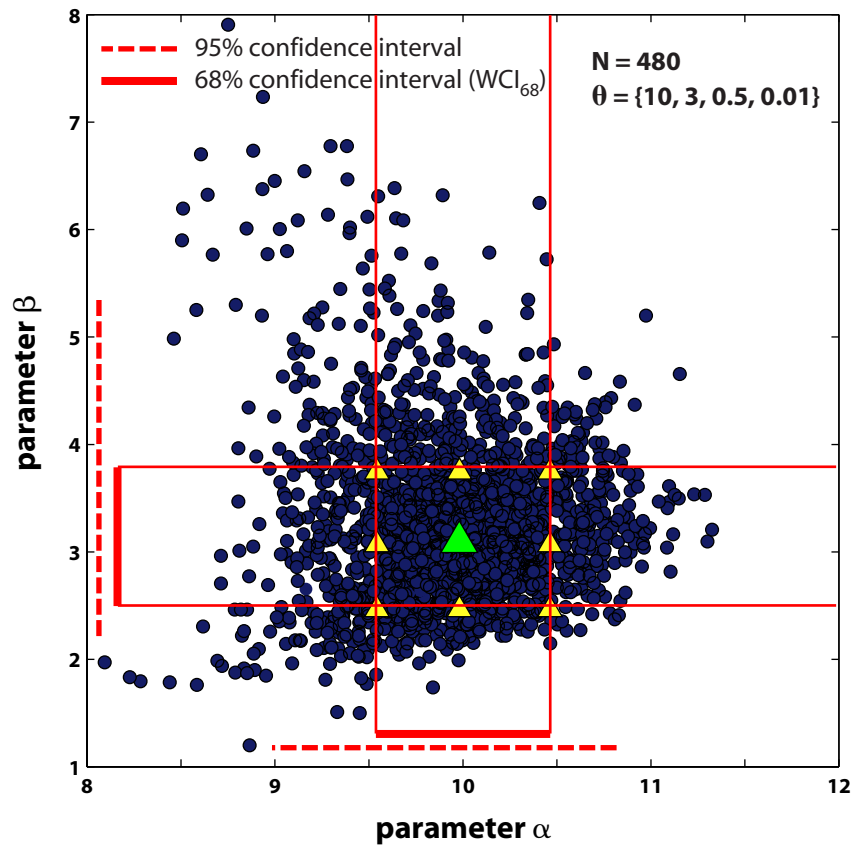


Figure 2:

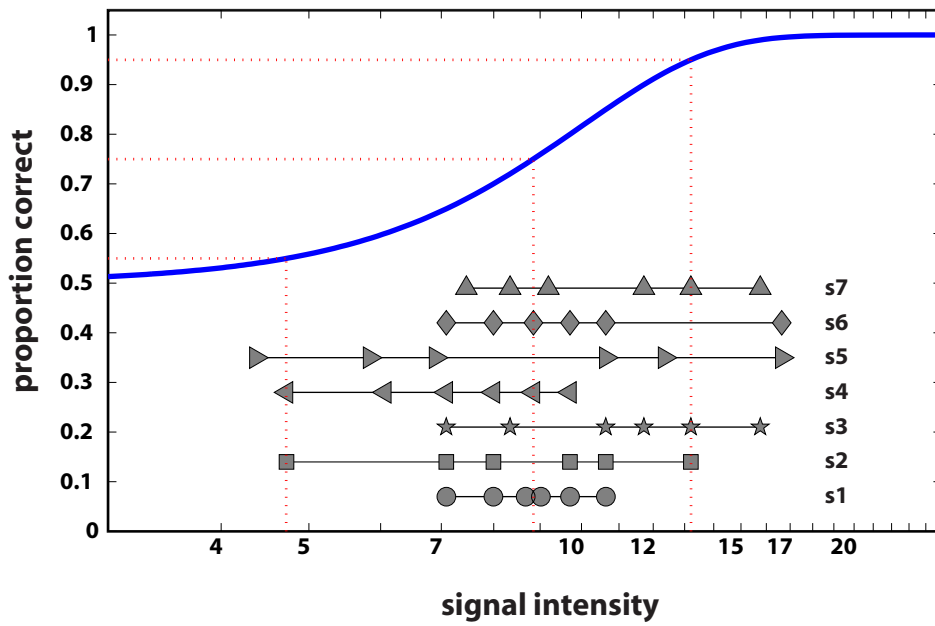


Figure 3:

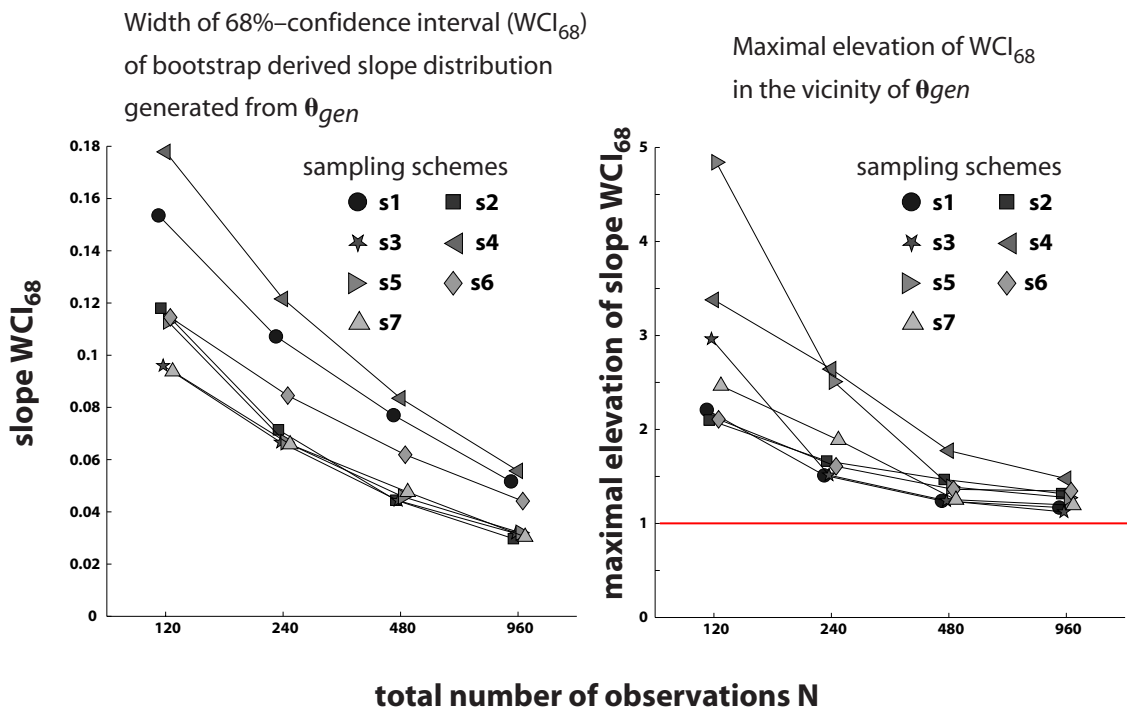


Figure 4:

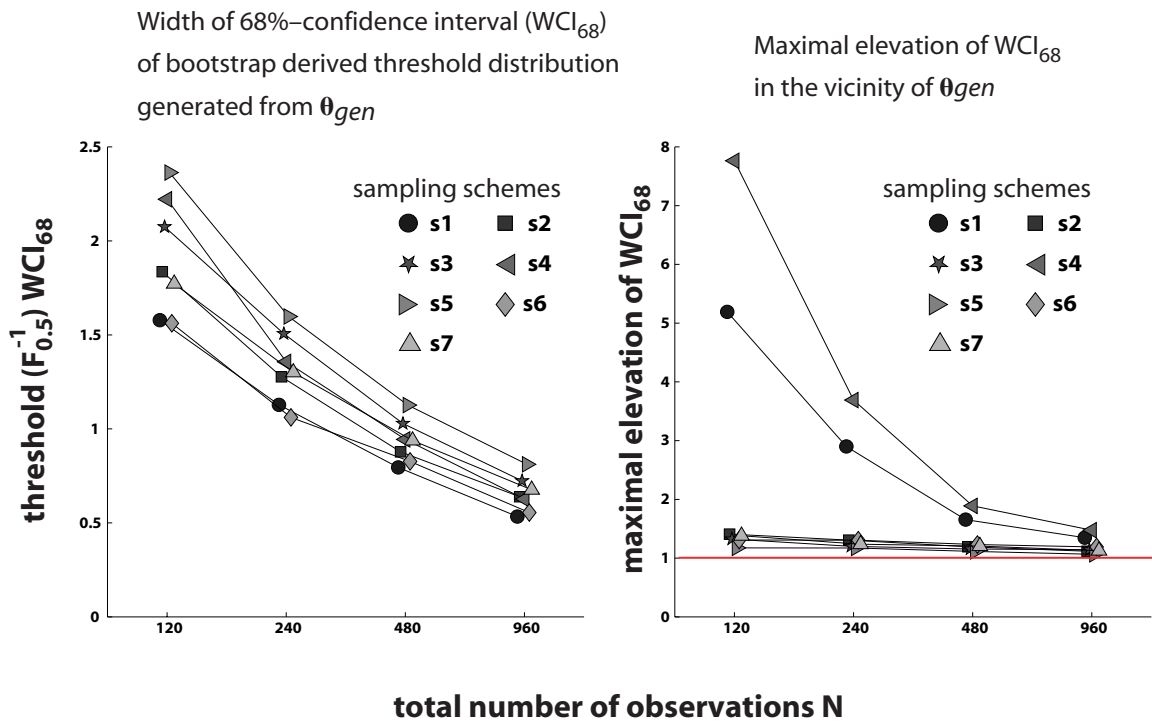


Figure 5:

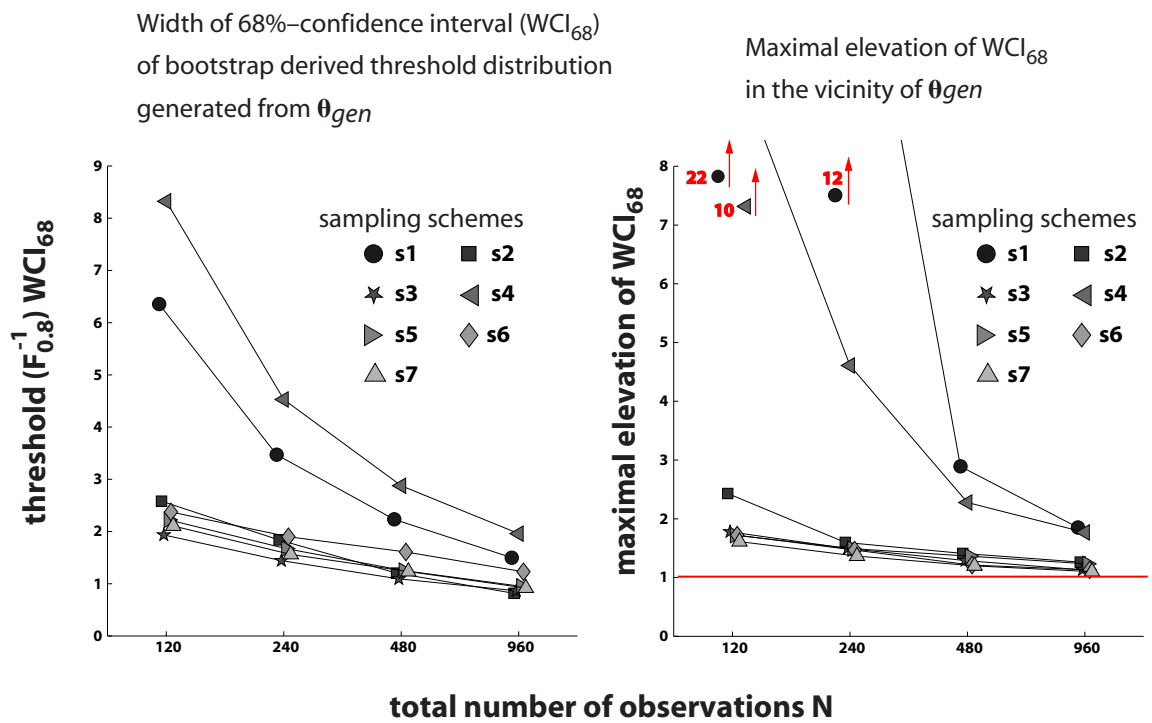


Figure 6:

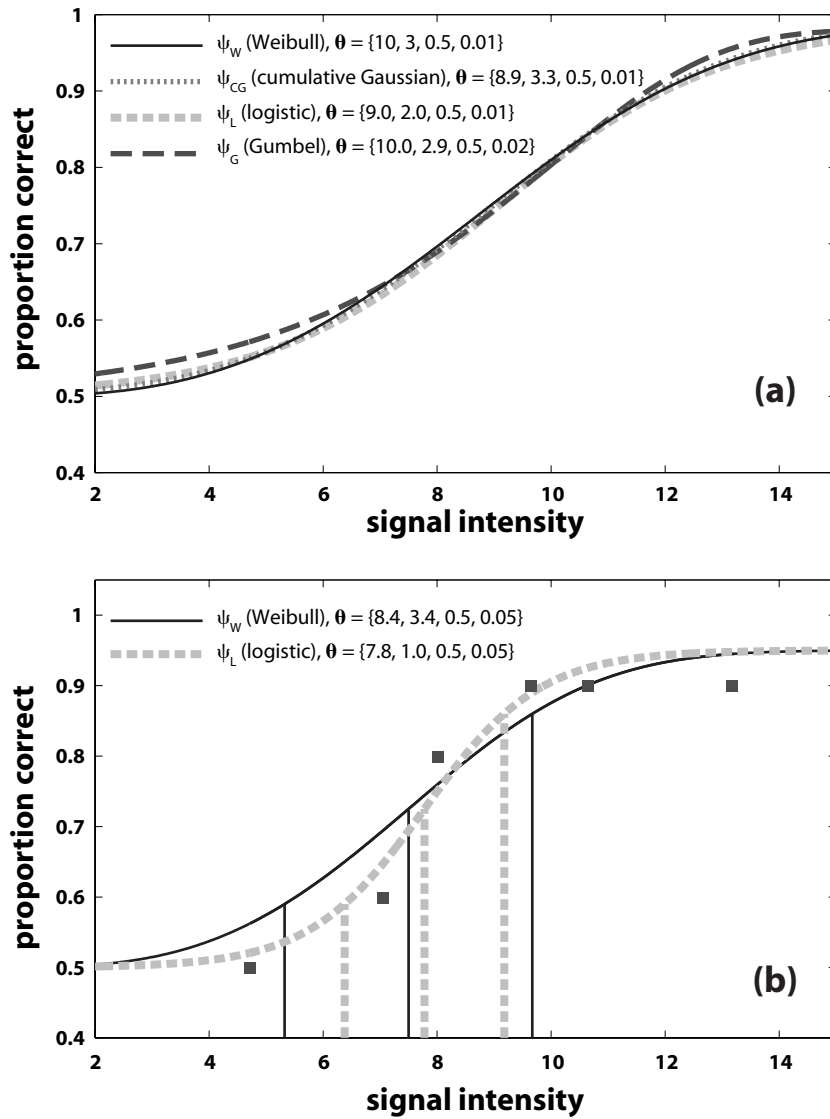


Figure 7

