

# Supplemental Materials for “Bayes Factors for Reinforcement-Learning Models of the Iowa Gambling Task”

Helen Steingroever<sup>a</sup>, Ruud Wetzels<sup>b</sup>, and Eric-Jan Wagenmakers<sup>a</sup>

<sup>a</sup> Department of Psychology, Psychological Methods, University of Amsterdam, The Netherlands

<sup>b</sup> PricewaterhouseCoopers, Amsterdam, The Netherlands

In these supplemental materials, we present a recipe on how to obtain Bayes factors with importance sampling, and two tests to check our implementation of importance sampling: (1) a model-recovery study, and (2) the Savage-Dickey density ratio test for each model. In addition, we present the results of a robustness analysis showing that our conclusions are unaffected by the choice of the priors on the model parameters. Finally, we present a model comparison study using BIC for the same models and data pool as used in the article.

## Recipe for Importance Sampling

In this section, we present a recipe that describes how we obtained Bayes factors with importance sampling. We use  $\mathcal{M}_{(\cdot)}$  to refer to specific model that can either be the EV, PVL, PVL-Delta, or the VPP model.

1. Fit model  $\mathcal{M}_{(\cdot)}$  to the data of participant  $s = 1$ .
2. Find the beta distributions (i.e.,  $\text{Beta}(\alpha, \beta)$ ) with the best fit to the posterior distributions of  $\theta$ .<sup>1</sup> Save the corresponding  $\alpha$  and  $\beta$  parameters.

---

<sup>1</sup>Note that  $\theta$  represents a subject and model-specific parameter vector (see Table 2 in the main article for each model’s parameters). This means that we obtain one beta distribution for each of the parameters contained in  $\theta$ .

3. Draw a set of parameters from the Beta importance densities, and compute the associated likelihood. Save the likelihood.
4. Repeat the previous step  $D - 1$  times (with  $D$  the number of draws).
5. Compute the marginal likelihood using Equation 1.

$$m(y | \mathcal{M}_{(\cdot)}) \approx \frac{1}{N} \sum_{i=1}^N \frac{p(y | \theta_i, \mathcal{M}_{(\cdot)})p(\theta_i | \mathcal{M}_{(\cdot)})}{g(\theta_i | \mathcal{M}_{(\cdot)})}, \quad \theta_i \sim g(\theta | \mathcal{M}_{(\cdot)}) \quad (1)$$

6. Repeat steps 1 – 5 for all  $s \in \{2, \dots, S\}$  (with  $S$  the number of participants).

### Model-Recovery Studies

In this section, we present the results of the model-recovery study. The purpose of this study was to confirm that the Bayes factor tends to favor the data-generating model. This study is based on eight generated data sets: We generated 25 synthetic participants completing a 100-trial IGT using each of the four models. As data-generating parameters we used the median parameter values obtained from fitting the models to a subset of the data used in the article.

We fit each of the four models to the four data sets, and then applied importance sampling to derive Bayes factors for all possible model comparisons. Analogous to the analyses reported in the main text, we present histograms showing the distribution of the log Bayes factors. In addition, we calculated the median posterior model probability for each model, and the proportion of participants for whom each model has the highest posterior model probability. The latter two should be high whenever the data-generating model is the same as the model that was used to fit the data (see also Pitt & Myung, 2002).

Figure 1 shows the distribution of the log Bayes factors of 25 synthetic participants completing a 100-trial IGT. It is evident that in the case of all models, the majority of the synthetic participants provides evidence for the data-generating model. This finding is corroborated by Table 1: The median posterior model probability and the percentage of participants for whom each model has the largest posterior model probability are highest

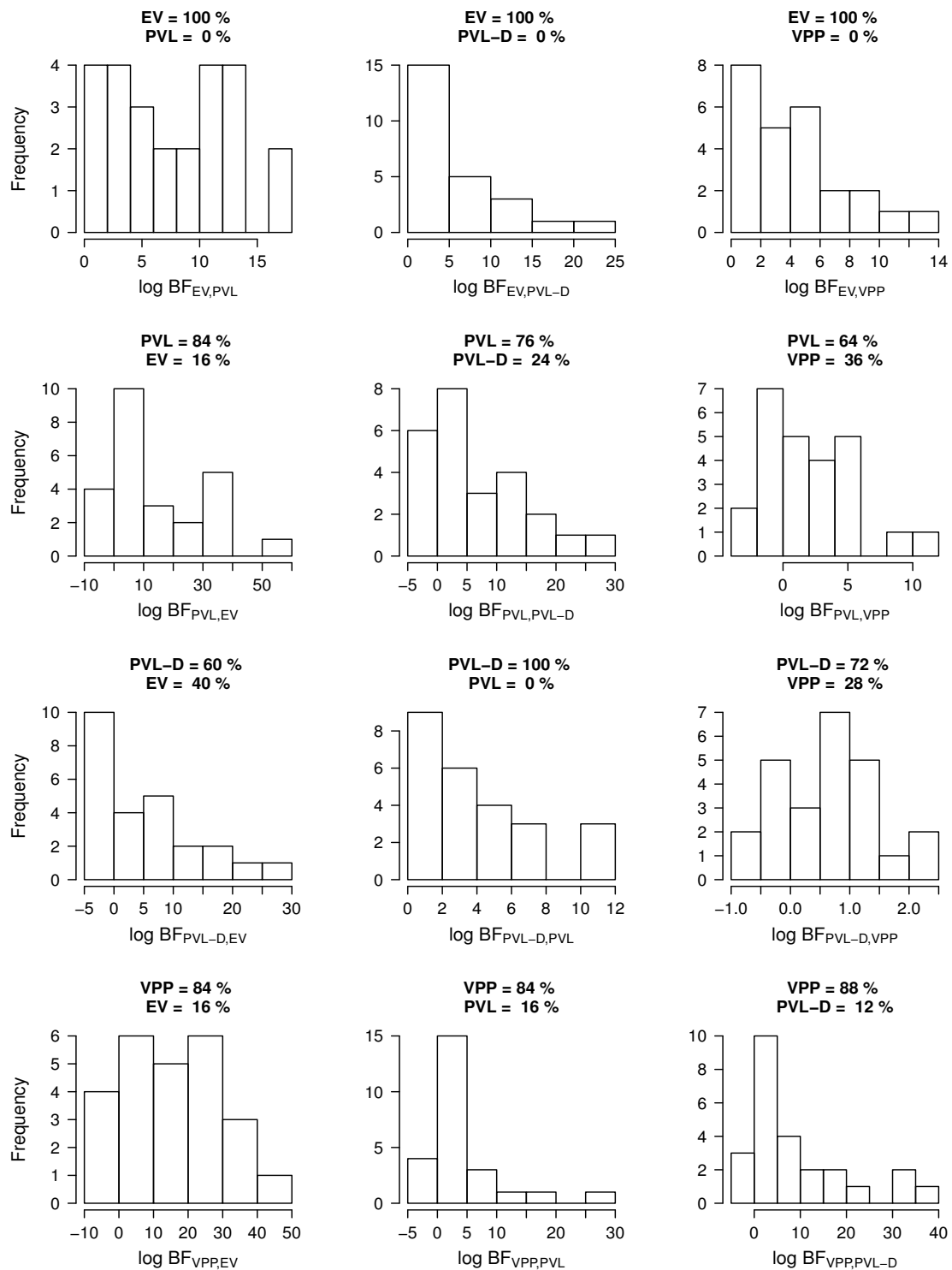


Figure 1. Histograms of the log(BF) of 25 synthetic participants completing a 100-trial IGT. Data of the first to fourth row were generated with the EV, PVL, PVL-Delta, and VPP model, respectively. A positive log(BF<sub>12</sub>) indicates that the data are more likely to occur under the first model (i.e., the data-generating model) than under the second model, whereas a negative log(BF<sub>12</sub>) indicates that the data are more likely to occur under the second model (i.e., the model that did *not* generate the data).

for the data-generating model. Thus, these results suggest that our implementation of importance sampling is correct and that the Bayes factor is a useful model comparison tool.

Table 1

*Median posterior model probabilities (MPMP) and percentage of participants for whom the corresponding model has the largest posterior model probability. The data were generated with either the EV, PVL, PVL-Delta, or VPP model, and describe the performance of 25 synthetic participants on a 100-trial IGT (i.e., first model-recovery study).*

	Data-generating model							
	EV		PVL		PVL-Delta		VPP	
	MPMP	%	MPMP	%	MPMP	%	MPMP	%
EV	.95	100	.00	12	.07	40	.00	16
PVL	.00	0	.69	56	.05	0	.07	12
PVL-Delta	.03	0	.02	12	.58	56	.02	8
VPP	.02	0	.12	20	.21	4	.69	64

### Savage-Dickey Density Ratio Tests

An alternative way to check our implementation of importance sampling is to investigate whether Bayes factors obtained with our implementation of importance sampling are in line with Bayes factors obtained with the Savage-Dickey density ratio test (Dickey, Lientz, et al., 1970; Dickey, 1971). The Savage-Dickey density ratio offers a method to compute Bayes factors for *nested* models. In order to be able to compare Bayes factors obtained with these two different methods, we thus needed to create nested RL models. This was done by fixing an arbitrary parameter of each model. We decided to fix the  $a$  parameter of each model to a predefined value  $a_0$ , and indicate nested models by  $\mathcal{M}_{(\cdot)}^*$ . Thus, the idea is to compare each of the four RL models to its nested version using both importance sampling and Savage-Dickey.

The Savage-Dickey method is explained in detail in Lee and Wagenmakers (2013), Vandekerckhove, Matzke, and Wagenmakers (2015), and Wagenmakers, Lodewyckx, Kuriyal, and Grasman (2010); here, we only provide the main idea: To obtain a Bayes factor comparing a RL model  $\mathcal{M}_i$  (where  $i \in \{\text{EV, PVL, PVL-Delta, VPP}\}$ ) to its nested version  $\mathcal{M}_i^*$ , we need to divide the prior ordinate at a fixed value of parameter  $a$  (i.e.,  $a_0$ ) by the posterior ordinate at that same fixed parameter value. The Bayes factor according

to the Savage-Dickey method is then defined as:

$$\text{BF}_{\mathcal{M}_i \mathcal{M}_i^*} = \frac{p(y | \mathcal{M}_i)}{p(y | \mathcal{M}_i^*)} = \frac{p(a = a_0 | \mathcal{M}_i)}{p(a = a_0 | y, \mathcal{M}_i)}, \quad (1)$$

where  $y$  is the data, and  $a = a_0$  indicates that the parameter  $a$  is fixed to a predefined value  $a_0$ .

The Bayes factor that we wish to approximate with importance sampling is the ratio of the marginal likelihood of the complete RL model and its nested version, that is:

$$\text{BF}_{\mathcal{M}_i \mathcal{M}_i^*} = \frac{m(y | \mathcal{M}_i)}{m(y | \mathcal{M}_i^*)}. \quad (2)$$

We applied the Savage-Dickey density ratio test and importance sampling to the same synthetic data set as used in the last section (i.e., 25 synthetic participants completing a 100-trial IGT). In Figure 2 we present the Savage-Dickey density ratio test for the four models and four synthetic subjects; the results for the remaining participants are similar. The header of each plot shows the Bayes factor obtained with importance sampling (i.e., BF\_IS), and the Bayes factor obtained with the Savage-Dickey method (i.e., BF\_SD). The dashed and solid lines represent the prior and posterior distribution, respectively. The black dots indicate the height of the prior and posterior distributions at  $a = a_0$ . From the figure it is evident that there is a close correspondence between Bayes factors obtained with the Savage-Dickey density ratio test and importance sampling suggesting that we correctly implemented importance sampling.

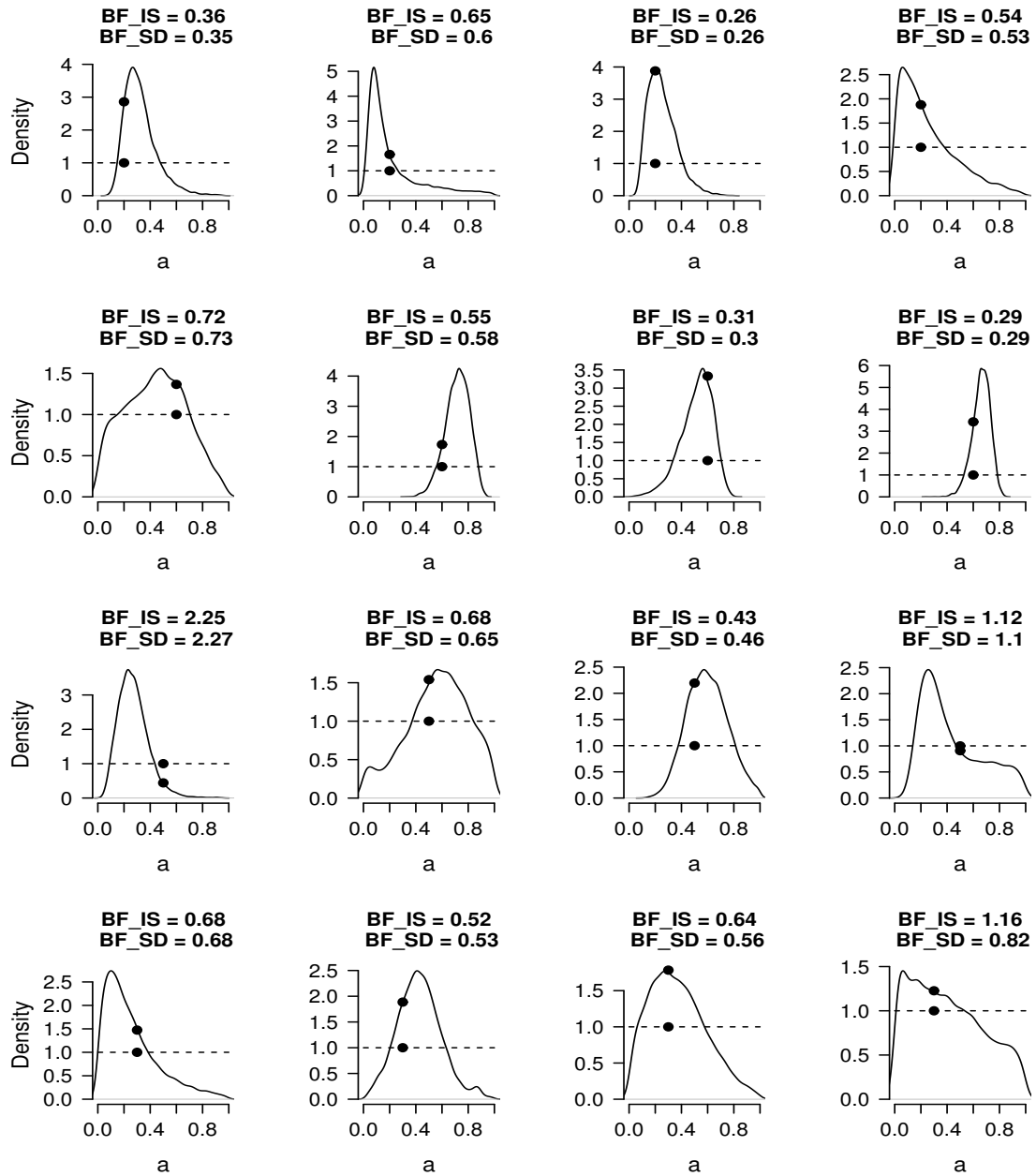
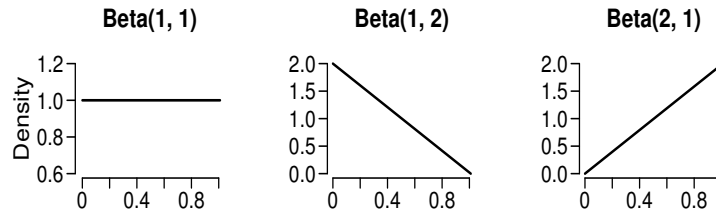


Figure 2. Illustration of the Savage-Dickey density ratio test for all models. Data of the first to fourth row were generated and fit with the EV, PVL, PVL-Delta, and VPP model, respectively. The header of each plot shows the BF obtained with importance sampling (i.e., BF\_IS), and the Bayes factor obtained with the Savage-Dickey method (i.e., BF\_SD). The dashed and solid lines represent the prior and posterior distribution, respectively. The black dots indicate the height of the prior and posterior distributions at  $a = a_0$ .

### Robustness Analyses

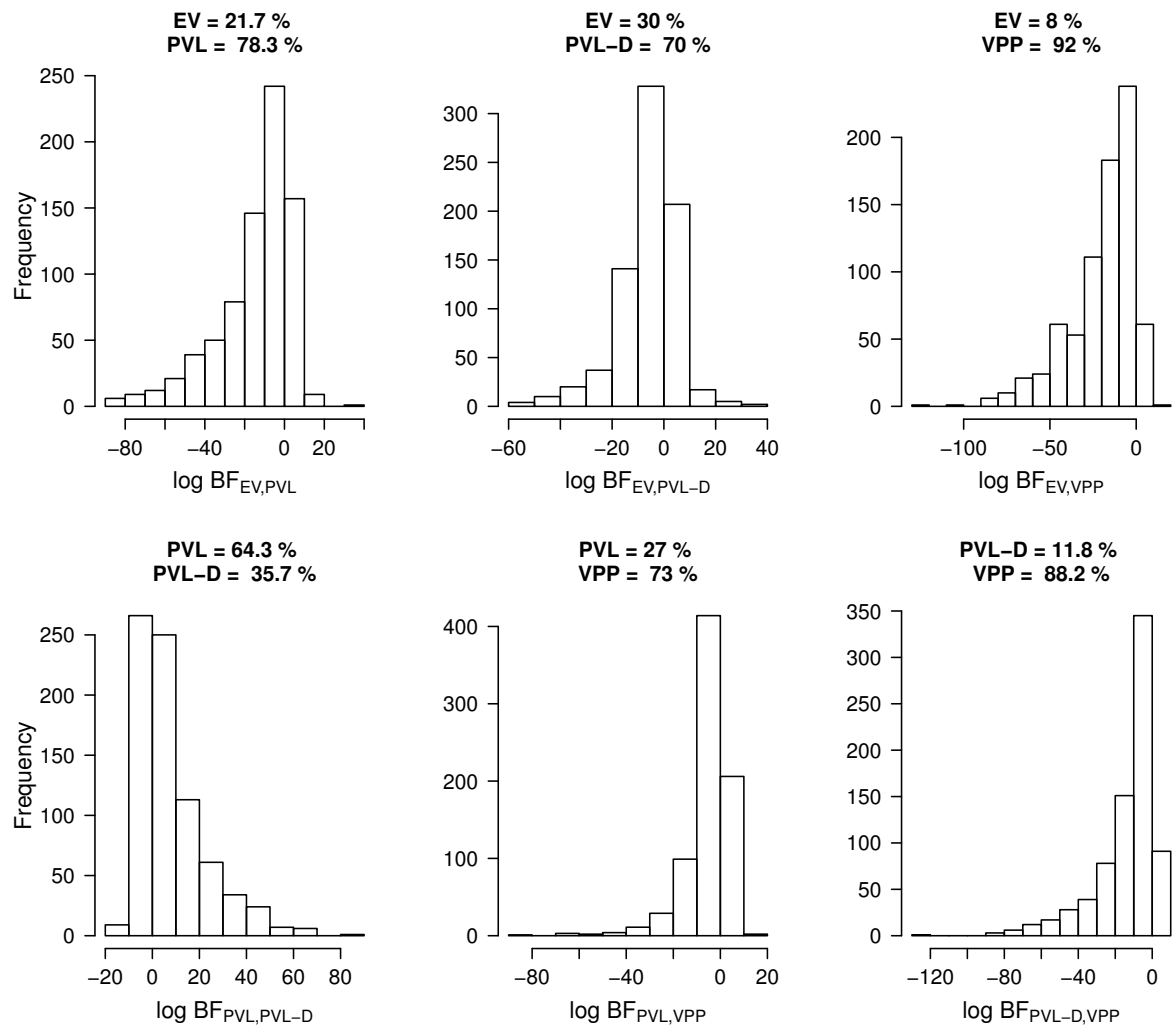
In this section, we present the results of a robustness analysis. The aim of this analysis is to investigate the extent to which our conclusions are altered by the choice of the priors on the model parameters. Whereas we used uniform priors on the model parameters in the analyses presented in the main article (i.e.,  $\text{Beta}(1, 1)$ ), we repeat here the analyses with two different priors: either a  $\text{Beta}(1, 2)$  or a  $\text{Beta}(2, 1)$  distribution. The different prior distributions are visualized in Figure 3. It is evident that the  $\text{Beta}(1, 1)$  distribution puts equal mass on all parameter values, the  $\text{Beta}(1, 2)$  distribution favors smaller parameter values, whereas the  $\text{Beta}(2, 1)$  distribution favors larger parameter values.



*Figure 3.* Visualization of the different priors. The prior distribution shown in the left panel is used in the analyses presented in the main article, whereas the prior distributions present in the middle and right panel are used in the sensitivity analyses.

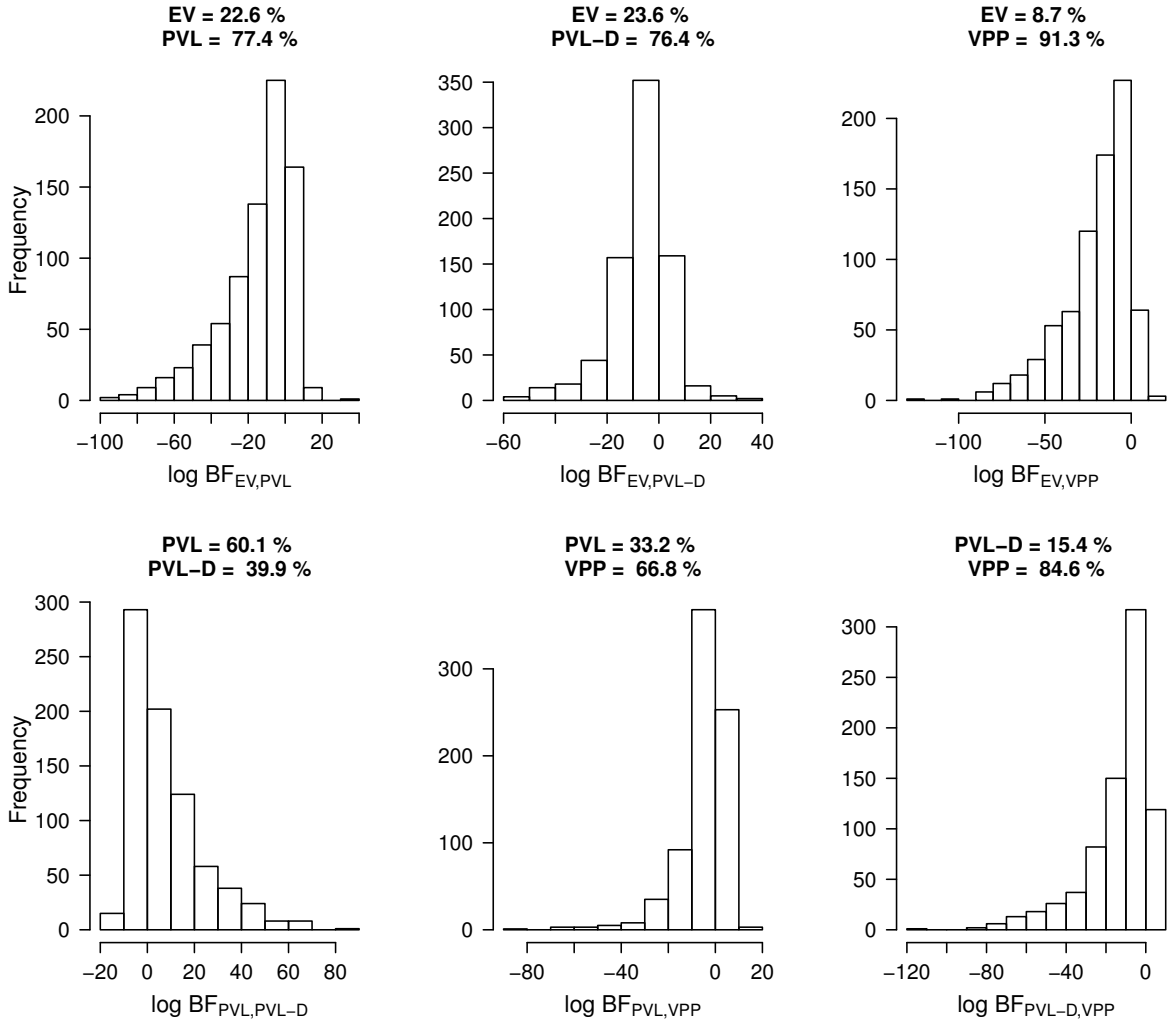
Figures 4 - 6 show, separately for the three different prior distributions, the distribution of the log Bayes factors of all participants for the six possible model comparisons. A positive  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the first model than under the second model, whereas a negative  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the second model than under the first model. The header of each histogram presents the percentage of participants for whom the data are more likely to occur under model  $\mathcal{M}_1$  than model  $\mathcal{M}_2$ .

Figures 4 - 6 show that there are some quantitative differences depending on which prior distribution is used. For example, the EV model is stronger supported when a  $\text{Beta}(2, 1)$  prior is used compared to the two other prior distributions. However, the qualitative

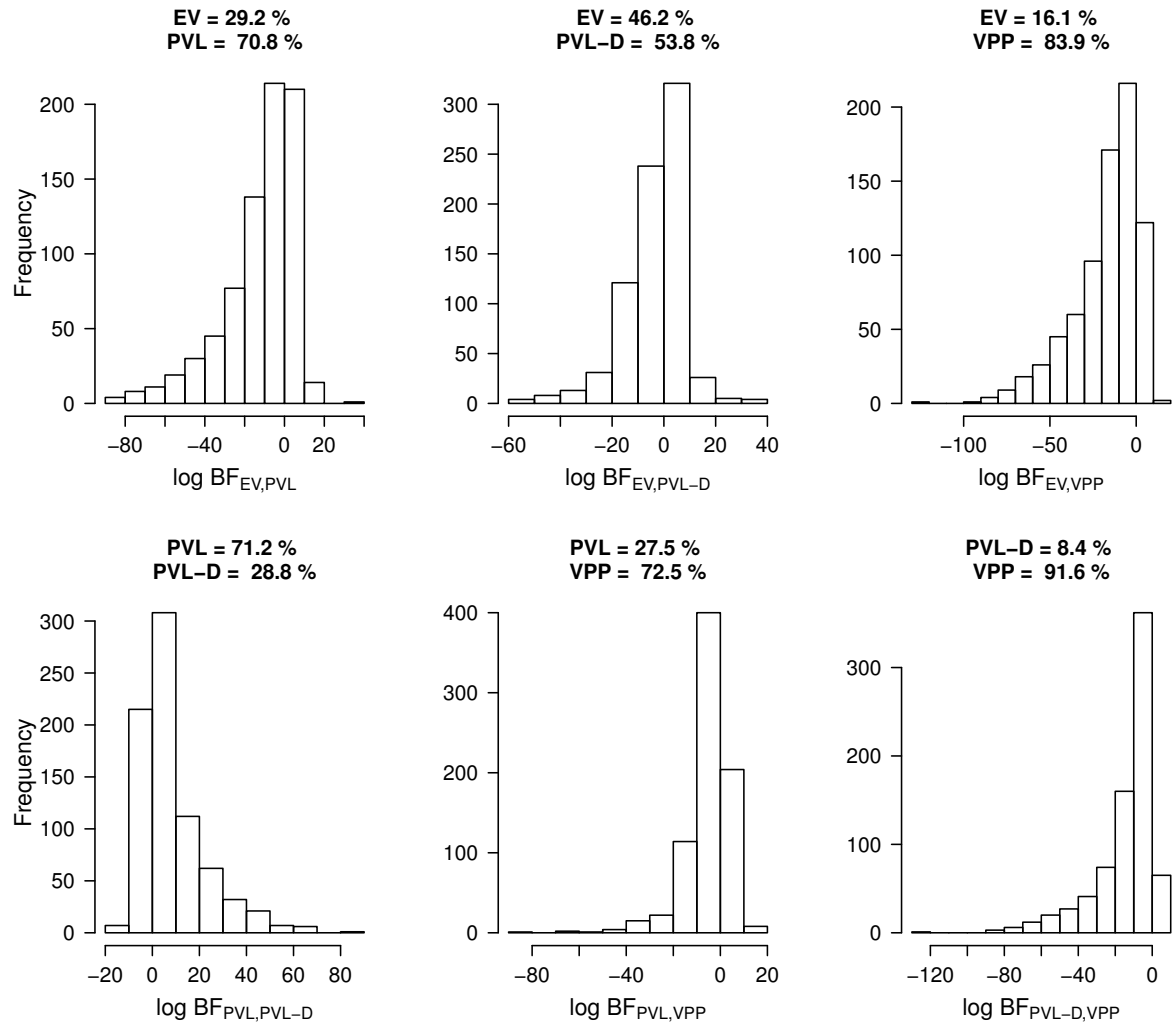


*Figure 4.* Beta(1, 1) prior: Histograms of the log(BF) for pairwise comparison of four RL models applied to the IGT data from each of 771 participants (cf. Figure 2 in the main article). A positive  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the first model than under the second model, whereas a negative  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the second model. Note that a  $\log(\text{BF})$  of 20 corresponds to a BF of almost 500 million, and that Jeffrey's (1961) considers as extreme evidence a Bayes factor larger than 100 (i.e.,  $\log(\text{BF}) > 4.6$ ). The header of each histogram presents the percentage of participants for whom the data are more likely to occur under the corresponding model.





*Figure 5.* Beta(1, 2) prior: Histograms of the  $\log(\text{BF})$  for pairwise comparison of four RL models applied to the IGT data from each of 771 participants. A positive  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the first model than under the second model, whereas a negative  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the second model. Note that a  $\log(\text{BF})$  of 20 corresponds to a BF of almost 500 million, and that Jeffreys (1961) considers as extreme evidence a Bayes factor larger than 100 (i.e.,  $\log(\text{BF}) > 4.6$ ). The header of each histogram presents the percentage of participants for whom the data are more likely to occur under the corresponding model.



*Figure 6.* Beta(2, 1) prior: Histograms of the  $\log(\text{BF})$  for pairwise comparison of four RL models applied to the IGT data from each of 771 participants. A positive  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the first model than under the second model, whereas a negative  $\log(\text{BF}_{12})$  indicates that the data are more likely to occur under the second model. Note that a  $\log(\text{BF})$  of 20 corresponds to a BF of almost 500 million, and that Jeffreys (1961) considers as extreme evidence a Bayes factor larger than 100 (i.e.,  $\log(\text{BF}) > 4.6$ ). The header of each histogram presents the percentage of participants for whom the data are more likely to occur under the corresponding model.

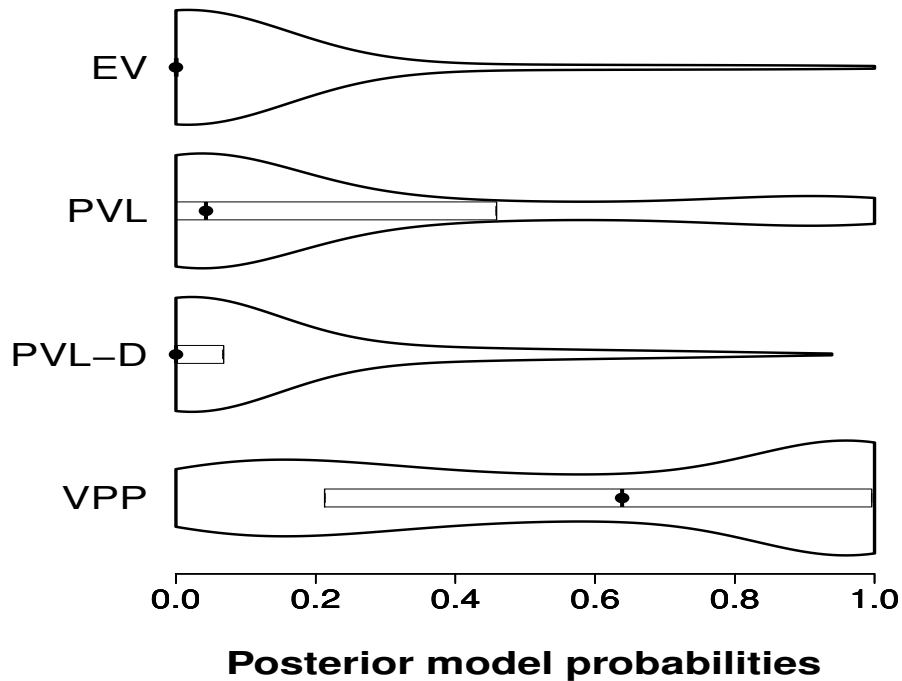
Table 2

*Median posterior model probabilities (MPMP; note that these need not sum to 1), and percentage of participants for whom the corresponding model has the largest posterior model probability, for the three different prior distributions separately. Grey shaded cells refer to the best model.*

	Beta(1, 1)		Beta(1, 2)		Beta(2, 1)	
	MPMP	%	MPMP	%	MPMP	%
EV	.00	7	.00	6	.00	13
PVL	.04	25	.04	30	.02	25
PVL-Delta	.00	9	.00	11	.00	5
VPP	.64	59	.49	52	.66	57

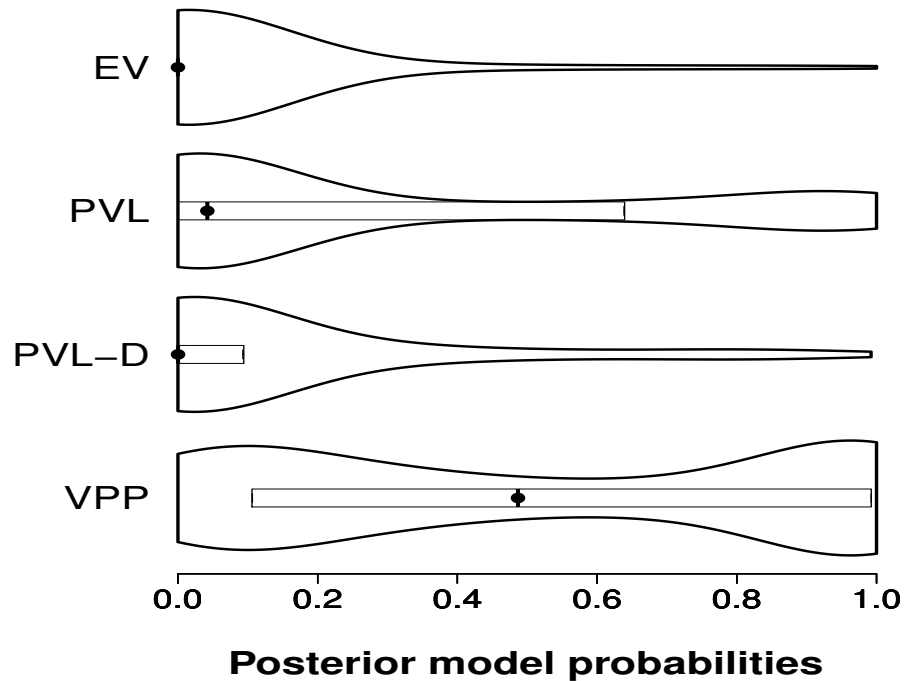
conclusions are the same irrespective of the prior distribution; all three figures show that the data provide the most evidence for the VPP model, and the least evidence for the EV model. In addition, the data provide more evidence for the PVL model than for the PVL-Delta model.

The findings from Figures 4 - 6 are corroborated by Table 2. The second, fourth, and sixth column of Table 2 show the median posterior model probabilities, and the third, fifth and seventh column show the percentage of participants for whom the corresponding model has the largest posterior model probability, separately for the three different prior distributions. It is evident that the VPP model is supported the most; that is, the data from 52-59% of the participants provide the most evidence for the VPP model. The PVL model is favored by the second largest proportion of the participants (i.e., 25-30%). It is also evident that the EV model is stronger supported than the PVL-Delta model when a Beta(2, 1) prior is used—a finding that is reversed for the two other prior distributions.



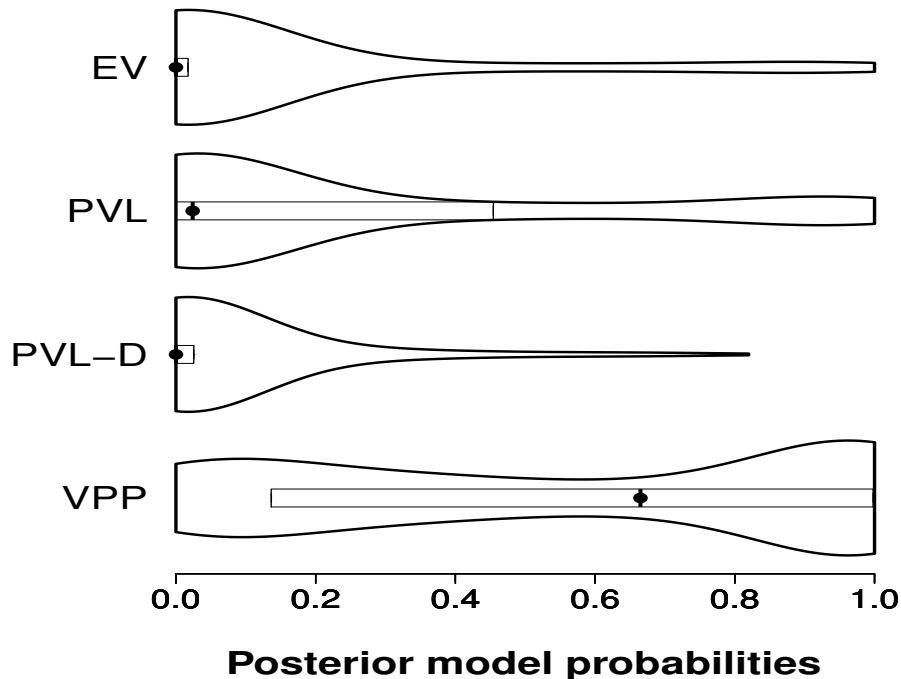
*Figure 7.* Beta(1, 1) prior: Distribution of the posterior model probabilities of 771 participants derived with importance sampling. Each violin plot shows the distribution of posterior model probabilities for one model. The dots indicate the median posterior model probability (cf. second column of Table 2), and the boxes indicate the interquartile range (i.e., the distance between the .25 and .75 quantiles).

The distributions of individual posterior model probabilities are visualized in Figures 7 - 9, which presents violin plots of the 771 posterior model probabilities for each of the four RL models, for the three different prior distributions separately. The dots indicate the median posterior model probability (cf. second, fourth, and sixth column of Table 2), and the boxes indicate the interquartile range (i.e., the distance between the .25 and .75 quantiles). From Figures 7 - 9, it is evident that in the case of the EV, PVL, and PVL-Delta models, the individual posterior model probabilities follow a right skewed distribution suggesting that the data of most participants provide little evidence for these models. It is also evident that the tail of the distribution in the case of the EV and PVL-Delta models is thinner than in the case of the PVL model. This suggests that there are more participants who provide strong



*Figure 8.* Beta(1, 2) prior: Distribution of the posterior model probabilities of 771 participants derived with importance sampling. Each violin plot shows the distribution of posterior model probabilities for one model. The dots indicate the median posterior model probability (cf. fourth column of Table 2), and the boxes indicate the interquartile range (i.e., the distance between the .25 and .75 quantiles).

evidence for the PVL model then for the EV and PVL-Delta models. In the case of the VPP model, the distribution of the posterior model probabilities is bimodal with the right mode being more pronounced than the left mode. This distribution suggests that the evidence for the VPP model differs greatly across participants, but that most participants provide compelling evidence in favor of the VPP model. Altogether Figures 7 - 9 suggest that there are only minor difference in the distributions of individual poster model probabilities. To conclude, this robustness analysis suggests that our main conclusions are unaffected by the choice of the prior distribution.



*Figure 9.* Beta(2, 1) prior: Distribution of the posterior model probabilities of 771 participants derived with importance sampling. Each violin plot shows the distribution of posterior model probabilities for one model. The dots indicate the median posterior model probability (cf. sixth column of Table 2), and the boxes indicate the interquartile range (i.e., the distance between the .25 and .75 quantiles).

### Comparison to BIC

In this section, we present the results of a model comparison study based on BIC for the same models and data pool as used in the article. The BIC is called post hoc fit criterion in the context of RL models for the IGT. Therefore, we call it here “BIC post hoc fit criterion”. The advantage of the BIC post hoc fit criterion is that it is easier to compute than the importance sampling Bayes factors. However, it should be kept in mind that the BIC post hoc fit criterion considers only one dimension of model complexity, that is, the number of parameters, and that the BIC post hoc fit criterion is derived as an asymptotic approximation of Bayesian model selection using Bayes factors (Myung, Cavagnaro, & Pitt, in press). Another popular measure is the Watanabe-Akaike

information criterion (WAIC; Watanabe, 2010, 2013). However, WAIC is not suitable for our predictive goal, that is, to predict the next choice given all previous choices (Aki Vehtari, personal communication, 16.07.2014; see also a discussion on Andrew Gelman’s blog <http://andrewgelman.com/2014/09/25/waic-time-series/>, and Vehtari & Ojanen, 2012).

**Computation of BIC.** The BIC for model  $M_{(\cdot)}$  is defined as follows (Schwarz, 1978):

$$\text{BIC}_{\mathcal{M}_{(\cdot)}} = -2 \log(L_{(\cdot)}) + k_i \log(n), \quad (3)$$

where  $L_{(\cdot)}$  is the maximum likelihood of model  $\mathcal{M}_{(\cdot)}$ ,  $k_{(\cdot)}$  is the number of free parameters of model  $\mathcal{M}_{(\cdot)}$ , and  $n$  is the number of IGT trials (Wagenmakers, 2007; Worthy, Pang, & Byrne, 2013, but see also for example Ahn, Bussemeyer, Wagenmakers, & Stout, 2008, Fridberg et al., 2010, and Yechiam, Arshavsky, Shamay-Tsoory, Yaniv, & Aharon, 2010, where the BIC post hoc fit criterion is computed for RL models relative to a baseline model). Thus, the first term in Equation 3 (i.e., the log maximum likelihood) quantifies the goodness-of-fit, whereas the second term penalizes a model for its complexity. Note that for the sake of clarity we omitted the notation that indexes a specific participant.<sup>2</sup>

**Approximation of the Bayes Factor.** The BIC score can be used to approximate the Bayes factor using the following equation (e.g., Wagenmakers, 2007):

$$\text{BF}_{12} \approx \exp\left(\frac{\text{BIC}_{\mathcal{M}_2} - \text{BIC}_{\mathcal{M}_1}}{2}\right), \quad (4)$$

Equation 4 allows us to investigate whether the approximations of the Bayes factors are in line with Bayes factors obtained from importance sampling.

---

<sup>2</sup>Since we did not use maximum likelihood to estimate the parameters, the fitting routine did not automatically provide us with  $L_{(\cdot)}$ —the maximum likelihood of model  $\mathcal{M}_{(\cdot)}$ . However, we obtained  $L_{(\cdot)}$  by computing the likelihood of the parameter combination that corresponds to the maximum log posterior. The log posterior is automatically returned by Stan (i.e., called “lp\_\_”). The BIC computation was confirmed by comparing our results obtained for the dataset of Worthy et al. (2013) to the ones reported in the original article.

Table 3

Median posterior model probabilities (MPMP), and percentage of participants for whom the corresponding model has the largest posterior model probability obtained from three different methods: (1) Importance sampling, and (2) BIC. Grey shaded cells refer to the best model.

	Importance Sampling		BIC	
	MPMP	%	MPMP	%
EV	.00	7	.00	14
PVL	.04	25	.36	46
PVL-Delta	.00	9	.00	17
VPP	.64	59	.00	24

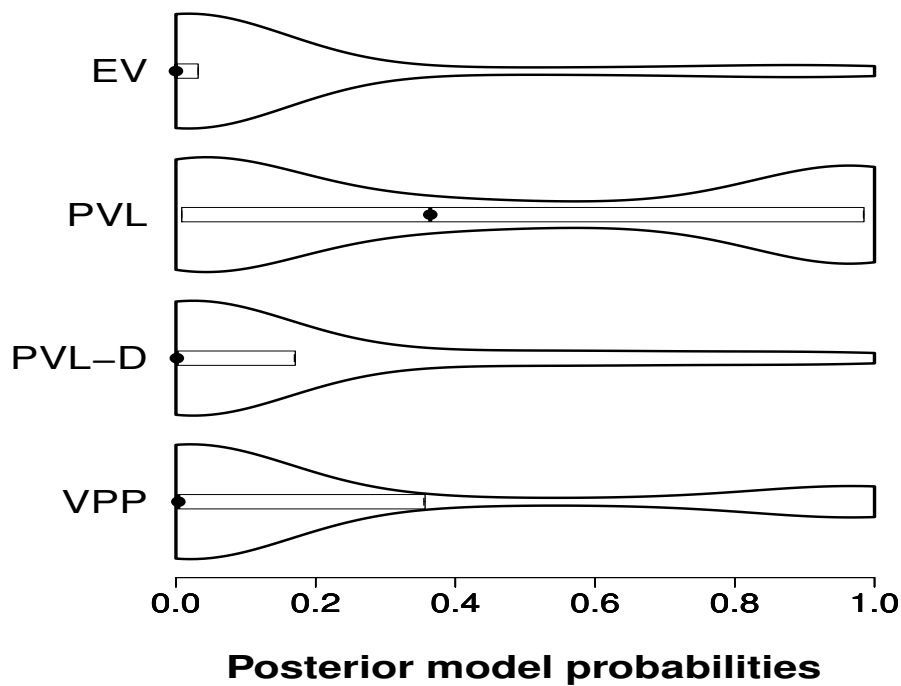


Figure 10. Distribution of the posterior model probabilities of 771 participants derived with BIC. Each violin plot shows the distribution of one model. The dots indicate the median posterior model probability (cf. Table 3), and the boxes indicate the interquartile range (i.e., the distance between the .25 and .75 quantiles).



**Results.** Table 3 shows the median posterior model probabilities (MPMP), and percentage of participants for whom the corresponding model has the largest posterior model probability obtained from two different methods: (1) Importance sampling, and (2) BIC post hoc fit criterion. Just as the Bayes factors obtained from importance sampling, the Bayes factors approximated with the BIC post hoc fit criterion suggest that the data of only a minority of participants provide strong evidence for the EV and PVL-Delta models. However, it is evident that in contrast to the Bayes factors obtained from importance sampling, Bayes factors approximated with the BIC post hoc fit criterion suggest that the data provide the most evidence for the PVL model and relatively little evidence for the VPP model. These findings are corroborated by Figure 10 showing the distributions of the posterior model probabilities of all participants derived with the BIC post hoc fit criterion. This analysis illustrates the critique that the BIC prefers simple models that underfit the data (Burnham & Anderson, 2002). In this particular case, the VPP model is punished for having relatively many parameters; however our Bayes factor analysis reveals that for this specific model comparison exercise, the number of free parameters alone is a limited and possibly misleading index of model complexity.

#### References

- Ahn, W.-Y., Busemeyer, J. R., Wagenmakers, E.-J., & Stout, J. C. (2008). Comparison of decision learning models using the generalization criterion method. *Cognitive Science*, *32*, 1376 - 1402.
- Burnham, K. P., & Anderson, D. R. (2002). *Model selection and multimodel inference: A practical information-theoretic approach*. Springer, New York.
- Dickey, J. M. (1971). The weighted likelihood ratio, linear hypotheses on normal location parameters. *The Annals of Mathematical Statistics*, 204 - 223.
- Dickey, J. M., Lientz, B., et al. (1970). The weighted likelihood ratio, sharp hypotheses about chances, the order of a Markov chain. *The Annals of Mathematical Statistics*, *41*, 214 - 226.
- Fridberg, D. J., Queller, S., Ahn, W.-Y., Kim, W., Bishara, A. J., Busemeyer, J. R., . . . Stout, J. C. (2010). Cognitive mechanisms underlying risky decision-making in chronic cannabis users. *Journal of Mathematical Psychology*, *54*, 28 - 38.
- Jeffreys, H. (1961). *Theory of probability* (Third ed.). Oxford University Press, Oxford, England.

- Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian modeling for cognitive science: A practical course*. Cambridge University Press Cambridge, MA.
- Myung, J. I., Cavagnaro, D. R., & Pitt, M. A. (in press). Model evaluation and selection. In J. Busemeyer, J. Townsend, Z. J. Wang, & A. Eidels (Eds.), *New Handbook of Mathematical Psychology, Vol. 1: Measurement and Methodology*. Cambridge, U.K.: Cambridge University Press.
- Pitt, M. A., & Myung, I. J. (2002). When a good fit can be bad. *Trends in Cognitive Sciences*, *6*, 421 - 425.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, *6*, 461 - 464.
- Vandekerckhove, J., Matzke, D., & Wagenmakers, E.-J. (2015). Model comparison and the principle of parsimony. In J. Busemeyer, J. Townsend, Z. J. Wang, & A. Eidels (Eds.), *Oxford Handbook of Computational and Mathematical Psychology*. Oxford: Oxford University Press.
- Vehtari, A., & Ojanen, J. (2012). A survey of Bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys*, *142* - 228.
- Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic Bulletin & Review*, *14*, 779 - 804.
- Wagenmakers, E.-J., Lodewyckx, T., Kuriyal, H., & Grasman, R. (2010). Bayesian hypothesis testing for psychologists: A tutorial on the Savage–Dickey method. *Cognitive Psychology*, *60*, 158 - 189.
- Watanabe, S. (2010). Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *The Journal of Machine Learning Research*, *3571* - 3594.
- Watanabe, S. (2013). A widely applicable Bayesian information criterion. *The Journal of Machine Learning Research*, *14*, 867 - 897.
- Worthy, D. A., Pang, B., & Byrne, K. A. (2013). Decomposing the roles of perseveration and expected value representation in models of the Iowa gambling task. *Frontiers in Psychology*, *4*.
- Yechiam, E., Arshavsky, O., Shamay-Tsoory, S. G., Yaniv, S., & Aharon, J. (2010). Adapted to explore: Reinforcement learning in Autistic Spectrum Conditions. *Brain and Cognition*, *72*, 317 - 324.