

Learning of Time Varying Functions is Based on Association Between Successive Stimuli

Lee-Xieng Yang (lxyang@nccu.edu.tw)

Department of Psychology, Researcher Center for Mind, Brain and Learning
National Chengchi University, No.64, Sec.2, ZhiNan Rd., Taipei City
11605, Taiwan (R.O.C)

Tzu-Hsi Lee (103752010@nccu.edu.tw)

Department of Psychology, National Chengchi University
No.64, Sec.2, ZhiNan Rd., Taipei City
11605, Taiwan (R.O.C)

Abstract

In function learning, the to-be-learned function is normally designed as time invariant. However, when the magnitudes of variable can be defined by time points, the function varies along time. Due to this difference in essence, the learning of the time-varying functions would be different from other functions. Specifically, the correlation between successive stimuli should play an important role for learning such functions. In this study, three experiments were conducted with the correlations set as positive high, negative high, and positive low. The results show people perform well when the correlation between successive stimuli is positive high or negative high. Also, people have difficulty learning the time-varying function with a low correlation between successive stimuli. A simple two-layered neural network model is evident to be able to provide good accounts for the data of all experiments. These results suggest that learning time varying function is based on association between successive stimuli.

Keywords: Function Learning; Time Varying Function

Function Learning

We are living in an orderly world, in which variables are mostly correlated with each other. For instance, the probability of rain might be a function of the extent to which the sky is overcast with dark clouds, or the distance to the car in front needed to avoid a car crash is a function of the current car speed. The study of how people learn a function and what people form to represent a learned function is referred to as function learning.

There are also two contrasting theoretical accounts in function learning. The rule-based account posits that people construct abstract rules to summarize the ensemble of experienced pairs of stimuli and responses used to teach the function. Most frequently, polynomial rules have been proposed as the representations of the mappings between stimulus magnitudes and response magnitudes (see Carroll, 1963; Koh & Meyer, 1991). On the contrary, the associative-based model assumes that people form direct associations between each stimulus and corresponding response without abstracting any summary information (Busemeyer, Byun, DeLosh, & McDaniel, 1997; DeLosh, Busemeyer, & McDaniel, 1997). However, the rule-based account overestimates the participants' performance in the extrapolation test but the associative-based model underestimates it. To get a better

theoretical account, a hybrid model combining these two approaches is proposed (McDaniel & Busemeyer, 2005).

Although these models differ on the assumption for the type of representation formed in function learning, it is basically agreed that the representation is formed for the whole function. However, contrary to this idea, it was found that people might form different representations for different parts of the function, such that a quadratic function was learned as the composition of two simpler monotonic functions, which were chosen for use at different contexts (Lewandowsky, Kalish, & Ngang, 2002). The POLE model (Kalish, Lewandowsky, & Kruschke, 2004) accounts for this finding well, by virtue of its architecture consisting of many modules, each of which represents a linear function corresponding only to a small region of the function, and a gating mechanism which always chooses one of the modules for use according to the stimulus value. Strictly speaking, the real function is not learned but approximated by the composition of many smaller linear functions.

Past studies have tested different functions and shown a number of characteristics of function learning. First, the linear functions are easier to learn than the nonlinear ones (see Busemeyer et al., 1997; Koh & Meyer, 1991). Second, it is found that it is more accurate to predict the response for the stimulus whose value falls in the training range (i.e., interpolation) than outside the range (i.e., extrapolation) (see Busemeyer et al., 1997; McDaniel & Busemeyer, 2005). Third, although the function of simpler forms (e.g., linear or power function) can be learned with the variables being of non-numeric forms (e.g., line length), Kalish (2013) reported that the periodic functions (e.g., sine function) cannot be learned without the employment of numeric stimuli. These characteristics reveal the limitations of human cognition for learning the functional relation between variables.

Time-Varying Function

Although many forms of functions have been tested, a particular form of function, which maps the timing of observation to the event at that timing seems not to have been tested yet. We call this function as time-varying function in this article, $y = f(t)$. An example of this function would be the height of water accumulated in a bucket from a constant supply source.

If the bucket is cylindrical, the height will be a linear function of time and if the bucket is conical, the height will be a parabolic function of time. To our knowledge, how people learn this kind of function has never been reported in literature. However, a relevant case in category learning has been reported recently.

Navarro and his colleagues tested how people could learn the categories when the category structure varies along training trials. In one of their experiments, the members of two categories moved up on the stimulus dimension constantly along with the increase of trial number and the categorization rule was set up as "Respond A, if $x_t > t$ and B otherwise" for any item x_t on trial t . Their results showed that participants could not only learn this category structure, but also be able to predict the item value on the next trial (Navarro & Perfors, 2009, 2012; Navarro, Perfors, & Vong, 2013). It is implied that people are able to capture some functional relationships between the time point (or trial number) and the stimulus value. However, the learning of the time-varying functions might be different from the normal functions.

Comparison Between Time-Varying Function and Normal Function

There are some features of the time-varying functions worth noting. First, due to that time can never return, when learning a time-varying function, making a prediction for response magnitude on each trial is always extrapolating what people have learned. However, in the case of learning the function $y = f(x)$, both the interpolation and extrapolation tests can be conducted.

Second, a time-varying function can be viewed as a function defining the relationships between successive stimuli, $x_t = f(x_{t-1})$. A good example is the game of throwing a Frisbee with friends. In this case, the only observable information is the spatial position of the Frisbee at any time point. Therefore, the best cue for us to estimate the position of the Frisbee at time t is its position at time $t - 1$.

Third, the learnability or complexity of function would be defined differently for the time-varying function. For the case of $y = f(x)$, the linear function has less parameters to estimate than the quadratic function, hence being easier to learn. For the case of $y = f(t)$, learning the functional relationship between time point to response magnitude is equivalent to learning to predict the next response magnitude with the current observed response magnitude. Thus, it is hypothesized that the time-varying function would be easy to learn, if the correlation between successive stimuli is high. If the correlation between successive stimuli is low, it would be hard to learn. To verify this hypothesis, three experiments were conducted.

Experiment 1

In this experiment, we first examined whether people can learn a linear time-varying function. The function was written as $x_t = t + \epsilon_t$, where t was trial number from 1 to 100 and ϵ

was randomly sampled from the uniform distribution between -0.5 and 0.5. All stimulus values were normalized between -15 and 15 for the convenience of computer programming. It was reasonable to expect that this function could be learned well, for (1) it was linear as well as (2) the correlation between successive stimuli was high.

Method

Participants and Apparatus There were in total 22 participants recruited from National Chengchi University in Taiwan for this experiment. Each participant was reimbursed by NTD\$ 60 (\approx US\$ 2) for their time and traffic expense. The whole experiment was conducted on an IBM compatible PC in a quiet booth. The processes of stimulus displaying and response recording were under the control of a computer script composed by PsychoPy (Peirce, 2007).

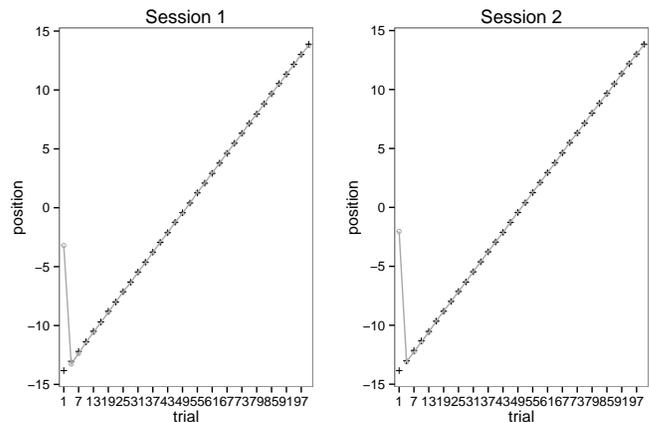


Figure 1: The stimulus structure in Experiment 1 (i.e., crosses) and the participants' predictions (i.e., circles) in Session 1 averaged across all participants.

Procedure The participants were instructed that they were playing a shooting game. In this game, they had to guess the position of a target on a horizontal line on the computer screen. On each trial, they moved the mouse cursor to where they thought the target would appear. After they pressed the space key to complete the guessing, the target would appear as an arrow on the correct position, together with a feedback text of "Hit" or "Miss" on the screen. The participants were told that "Hit" meant that your guess was close enough to the true answer and otherwise you would get "Miss". The whole experiment was conducted in two sessions, each of which consisted of 100 trials. The same 100 stimuli were presented in the two sessions. The distance between the target's correct position and the participants' guess was error. The amount of squared error and the proportion of received "Hit" (e.g., accuracy) were the dependent variable in this experiment.

Results

Visual inspection on Figure 1 shows that participants performed quite well except for the very early trials¹. For simplifying the complexity of data analysis, we divided the 100 stimuli to 10 blocks. The squared prediction error decreases from 40.29 to 0.03 with the mean = 4.06 through 10 blocks across two sessions. A Block (10) \times Session (2) within-subjects ANOVA reveals a significant main effect of Block on the squared error [$F(9, 189) = 72.83$, $MSe = 98$, $p < .01$], no significant main effect of Session [$F(1, 21) = 2.367$, $MSe = 166.30$, $p = .139$], and a significant interaction effect between Block and Session [$F(9, 189) = 2.346$, $MSe = 166.3$, $p < .05$].

The participant's accuracy is another dependent variable, which is computed as the number of "Hit" divided by all trials. Due to the "Hit" range was very small in our experiments, the highest accuracy in a block was .63 and the lowest was .36 across all sessions. A Block (10) \times Session (2) within-subjects ANOVA shows a significant main effect of Block on the accuracy [$F(9, 189) = 8.281$, $MSe = 0.028$, $p < .01$], no significant main effect of Session [$F(1, 21) < 1$], and a significant interaction effect between Block and Session [$F(9, 189) = 5.052$, $MSe = 0.027$, $p < .01$].

We also check the correlation between each participant's predictions and the true answers. The averaged Pearson's r across all participants is quite high [$r = .97$]. Together with the visual inspection on Figure 1, it is confirmed that people can learn the linear time-varying function very well.

Experiment 2

In this experiment, the function was set up as $x_t = 50 + (-1)^t \sqrt{100 - t}$, which made the target jump left and right, gradually moving toward the central point. Obviously, this function was far more complex than the one used in Experiment 1 and it was nonlinear. If the learning of $y = f(t)$ shared the same characteristics of the learning of $y = f(x)$, it should be expected that this function could not be learned well. However, if our discussion about the characteristics of time-varying function was right, it should be expected that this function could be learned well, due to high correlation between successive stimuli [$r = -.99$].

Method

Participants and Apparatus There were in total 21 participants recruited from National Chengchi University in Taiwan for this experiment. Each participant was reimbursed by NTD\$ 60 (\approx US\$ 2) for their time and traffic expense. The testing materials and procedure are all the same as those in Experiment 1.

¹For making the figure easier to read, we plot the human prediction by circles and the correct answers by crosses on only the even-numbered trials in the first session. The result pattern is the same in the second session.

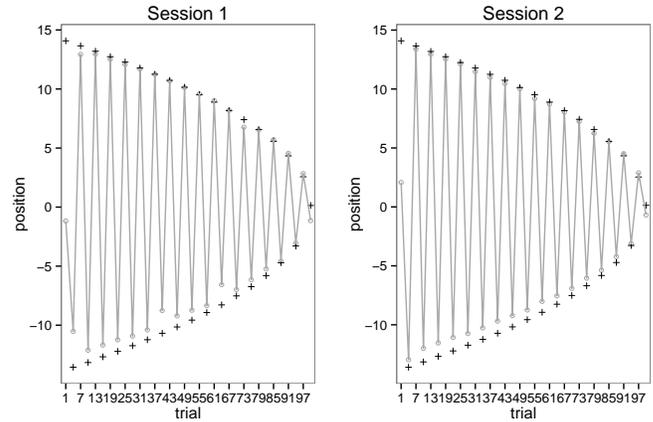


Figure 2: The stimulus structure in Experiment 2 (i.e., crosses) and the participants' predictions (i.e., circles) in Session 1 averaged across all participants.

Results

See the circles and crosses in Figure 2. Apparently, the participants could capture the moving pattern of the target, although on the early trials, they made some larger errors. Similar to what we found in Experiment 1, the squared prediction error drops along blocks from 73.79 to 1.57 (mean = 15.35) across two sessions. A Block (10) \times Session (2) within-subjects ANOVA reveals a significant main effect of Block [$F(9, 180) = 14.24$, $MSe = 1303$, $p < .01$], a significant main effect of Session [$F(1, 20) = 17.22$, $MSe = 196$, $p < .01$], and a significant interaction effect between Block and Session [$F(9, 180) = 16.12$, $MSe = 177.8$, $p < .01$]. Although the error curve goes down toward 0, the mean squared prediction error is 15.53 far larger than that in Experiment 1, which is 4.06. This suggests that the linear function is easier to learn than the quadratic function.

The accuracy data also suggest that this function is harder to learn than the linear function with the mean highest accuracy in a block across all participants and sessions as .34 and the lowest as .14. A Block (10) \times Session (2) within-subjects ANOVA reveals a significant main effect of block [$F(9, 180) = 9.747$, $MSe = 0.018$, $p < .01$], no significant main effect of Session [$F(1, 20) < 1$], and no significant interaction effect between Block and Session [$F(9, 180) < 1$].

Although the accuracy is quite low, this does not mean that people cannot learn this function. As shown in Figure 2, the participants' predictions are close to the true answers. Also, the correlation between each participant's predictions and the true answers is considerably high [mean $r = .92$]. As expected, the participants can learn this complex time-varying function.

Experiment 3

In this experiment, we would like to examine whether people could predict the stimulus magnitudes, when the corre-

lation between successive stimuli was lower. See Figure 3 as an example, which was the real case for testing one participant². The dashed line showed the true moving pattern of the stimulus, which was generated by $y = g[a] + z[b + 1]$, where $a = \lfloor ((t + 4)/5) \rfloor$, $b = t \bmod 5$, g was the random permutation of the vector [1,6,11,...,96], and for each g , z was a new random permutation of the vector [1,2,3,4,5]. The correlations between successive stimuli were averaged across all participants and all sessions as $r = .80$, which was lower than the correlations in the previous experiments. With no matter which view to look at this form (i.e., number of parameters to estimate or correlation between successive stimuli), it was expected that this function could not be learned well.

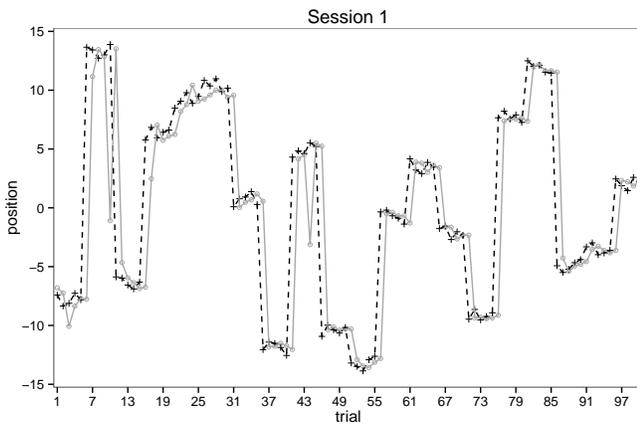


Figure 3: The stimulus structure in Experiment 3 (i.e., crosses) and predictions of participant #14 (i.e., circles).

Method

Participants and Apparatus There were in total 18 participants recruited for this experiment from National Chengchi University in Taiwan. Each participant was reimbursed by NTDS\$ 60 (\approx US\$ 2) for their time and traffic expense. The testing materials and procedure are all the same as those in Experiment 1.

Results

As shown in Figure 3, apparently, the participant could not predict the target position. Otherwise, we will see the dashed line (for answers) and solid line (for participant’s predictions) superimpose on each other. However, the response pattern is not random either. In fact, the participant’s predictions seem always to be one step behind the true answers. Although we do not show the predictions of the rest 17 participants, their predictions are one step behind the true answers also. Thus, strictly speaking, we do not think that the participants learned this function.

²Different participants received different moving patterns to learn.

The squared prediction error drops from 69.69 to 42.47 along blocks in Session 1 and has no clear change from 23.12 to 24.30 in Session 2. Although the performance gets better in Session 2, the prediction error never goes close to 0. The mean squared error for all participants across blocks and sessions is 30.844, which is larger than 15.53 (mean error in Experiment 2) and 4.06 (mean error in Experiment 1). Thus, the learning performance in this experiment is the worst among the three experiments in this study.

As done for the previous experiments, a Block (10) \times Session (2) within-subjects ANOVA was conducted for the prediction error. The results show no significant main effect of Block [$F(9, 153) = 1.53$, $MSe = 998.4$, $p = .142$], a significant main effect of Session [$F(1, 17) = 14.94$, $MSe = 424$, $p < .01$], and a significant interaction effect between Block and Session [$F(9, 153) = 3.206$, $MSe = 701.6$, $p < .01$].

The mean accuracy in a block across all sessions is even lower than that in the other two experiments. The highest mean accuracy is about .11 and the lowest is .06. It is clear that the participants cannot capture the moving pattern of the stimulus. A Block (10) \times Session (2) within-subject ANOVA shows no main effect of Block on accuracy [$F(9, 153) = 1.179$, $MSe = 0.006$, $p = .312$], no main effect of Session [$F(1, 17) = 3.367$, $MSe = 0.006$, $p = .08$], and no interaction effect between Block and Session [$F(9, 153) < 1$].

We also computed the Person’s r for each participant’s prediction and the true answer. Although the mean correlation is not low ($r = .76$), this finding might result from the fact that the participants’ prediction is always one step behind the true answer. To sum up, the linear function is the easiest to learn and the quadratic function is the second. Basically, participants cannot learn the complex function in Experiment 3. In order to get a better understanding about the underlying mechanism for learning the time-varying functions, we developed a neural network model for the learning of time-varying functions.

Model for Learning Time Varying Function

A time-varying function can be rewritten as $x_t = f(x_{t-1})$ and the simplest form of it would be $x_t = \beta_0 + \beta_1 x_{t-1}$. Thus, learning a time-varying function is equivalent to estimating the optimal parameter values, with which the model makes the smallest error. To this end, a simple two-layered neural network is proposed. There are two input nodes, which respectively correspond to the position of the stimulus on the preceding trial x_{t-1} and the standard moving distance which is set as 1. There is only one output node corresponding to the predicted position on the current trail $\hat{x}_t = w_1 \times 1 + w_2 x_{t-1}$. The associative weight w_1 represents the size of moving distance. The weight w_2 represents how much correlated the last position is with the current position. When the true answer x_t is provided, the error is then computed as $x_t - \hat{x}_t$.

The associative weights are updated with WH algorithm³

³This algorithm is a special case of backpropagation algorithm, which is specifically used for two-layered neural network models.

(Abdi, Valentin, & Edelman, 1999) to decrease the error made by the model. Also, we make the updating amount for weights decay all the way through training trials. Thus, the updated amount for w_1 on trial t is $\Delta w_{1,t} = \eta \exp^{-\xi(t-1)}(x_t - \hat{x}_t)$, where $\eta \geq 0$ is the learning rate and $\xi \geq 0$ determines how quickly the updated amount of weight drops. Likewise, $\Delta w_{2,t} = \eta \exp^{-\xi(t-1)}(x_t - \hat{x}_t)x_{t-1}$.

There are some features of this model worth noting. First, the associative weight w_2 actually reflects the correlation between successive stimuli. Second, this model only learns the correlation between successive stimuli and contains no summary information of the whole function. In fact, it can be applied to account for the learning of different time-varying functions, as no matter which form (complex or simple) the function has, the learning of a time-varying function can always be viewed as the learning of the association between successive stimuli. Thus, our model should be regarded as an associative-based model, not a rule-based model.

Modeling

The model was fit to each participant's data in each experiment with the stimulus positions being normalized between 0 and 1. Each participant's first response in each session was by default the first input for the model. The initial weights of w_1 and w_2 were set as 0 for all experiments except Experiment 3. The model provided the best fit for Experiment 3 data when w_2 was initially set as 1, suggesting that participants in Experiment 3 were more likely to repeat the observed position of stimulus on the preceding trail as the response for current trail. The statistics of optimally estimated parameter values and the goodness of fit (RMSD) for all experiments are listed in Table 1.

Table 1: Mean goodness of fit and mean estimated parameter values for a best fit with the standard deviation listed in parenthesis.

	RMSD	η	ξ
Exp 1	0.04 (0.02)	1.06 (0.71)	0.02 (0.09)
Exp 2	0.08 (0.03)	1.73 (1.14)	0.30 (0.55)
Exp 3	0.09 (0.03)	0.43 (0.55)	1.81 (4.14)

The smaller the RMSD, the better the fit is. Apparently, the model fit all the data very well. See the crosses in Figure 4, Figure 5, and Figure 6 for the model prediction in Session 1⁴, which are quite close to the circles denoting the participants' responses.

The estimated learning rate for Experiment 1 is about 1 and the decay rate is quite small, suggesting that decay of learning is not fast and leaning continues through training trials. The learned associative weights for the moving size $w_1 = 0.30$ and the correlation with the preceding stimulus $w_2 = 0.70$ suggest that the participants predict the current position of the target

⁴The pattern is almost the same for Session 2.

by moving it a certain distance (i.e., 0.30 times of the standard moving size) from the place a bit behind (i.e., 70%) the position just seen in the same direction of the last move.

For Experiment 2, the mean learning rate is high and so is the mean decay rate. This suggests that the model adjusts the associative weights largely on the early learning trials, but quickly halts doing so. The learned associative weights are $w_1 = 1.00$ and $w_2 = -0.94$. The negative weighting for the preceding position enables the model to make symmetrical predictions between successive trials and $|w_2| \leq 1$ enables the model to gradually converge the predicted position toward the midpoint.

For Experiment 3, the mean estimated learning rate is low and the decay rate is high, suggesting that the model has not updated the associative weights too much since early trials. In fact, the learned associative weights, $w_1 = 0.01$ and $w_2 = 0.98$, together suggest that the model merely repeats the preceding target position as the current prediction. As the model captures the participants' response patterns very well, it is implied that the participants did not actually learn the function but just repeated what they saw as the prediction for the next trial.

It is revealed in Experiment 2 that the larger η or ξ is, the smaller the error is ($r = -.51, p < .05$ for η and $r = -.57, p < .01$ for ξ) but no significant correlations between parameters and human performance in other experiments. This might be because that Experiment 1 and Experiment 3 are either too easy or too hard for the participants to learn.

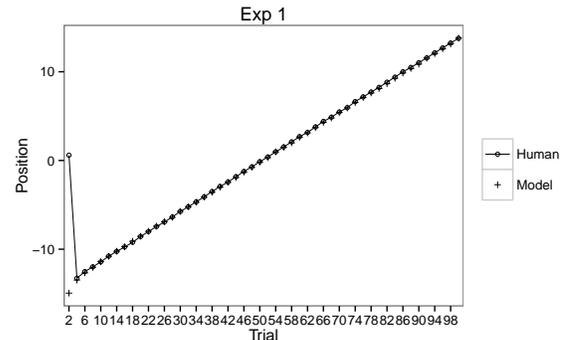


Figure 4: The model prediction and averaged human response in Session 1 in Experiment 1.

General Discussion

The main purpose of this study is to examine the characteristics of function learning with time-varying functions. Three experiments were conducted with different time-varying functions: linear, quadratic, and irregular. The differences between these functions are not only the complexity of the function form, but also the strength of correlation between successive stimuli. In the first two experiments, the correlation is very high regardless of the direction, whereas in the third experiment, the correlation is lower.

The behavioral data show that the learning of the linear and quadratic functions are easier than that of the irregular function, suggesting that the correlation between successive stimuli is critical to function learning with time-varying functions, not the number of parameters (or the complexity) of the function. The success of our model supports the associative-based account and implies that a time-varying function can be learned as a composition of many partial representations, not a holistic representation.

One may regard the learning of time-varying functions as operant conditioning. That may or may not be true, depending on what we think is actually conditioned. If the response is the target for conditioning, then the learning of time-varying functions is not operant conditioning, as every single response is new and it is impossible to reinforce the likelihood for the same response to be made in the future. However, if the moving size is the target for conditioning, then for the case in which the target moves constantly (e.g., the linear function in Experiment 1), we may regard the learning of the time-varying function as a kind of operant conditioning. However, for the case where the target moves in a decreasing (or increasing) speed (e.g., the quadratic function in Experiment 2), it might not be suitable to equate the learning of time-varying functions and operant conditioning. Future studies including the transfer trials are needed in order to examine whether people form any concept for the time-varying function.

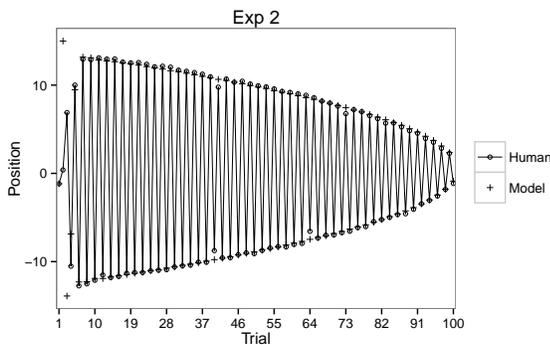


Figure 5: The model prediction and averaged human response in Session 1 in Experiment 2.

References

Abdi, H., Valentin, D., & Edelman, B. (1999). *Neural networks*. SAGE Publications, Inc.

Busemeyer, J. R., Byun, E., DeLosh, E., & McDaniel, M. A. (1997). *Learning functional relations based on experience with input-output pairs by humans and artificial neural networks* (K. Lamberts & D. R. Shanks, Eds.). Cambridge, MA, US: The MIT Press.

Carroll, J. D. (1963). *Function learning: The learning of continuous functional maps relating stimulus and response continua*. Princeton, NJ: Educational Testing Service.

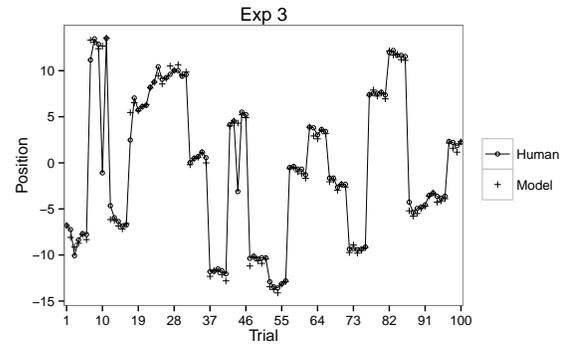


Figure 6: The model prediction and human response of participants #14 in Session 1 in Experiment 3.

DeLosh, E. L., Busemeyer, J. R., & McDaniel, M. A. (1997). Extrapolation: The sine qua non for abstraction in function learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *23*, 968-986.

Kalish, M. (2013). Learning and extrapolating a periodic function. *Memory & Cognition*, *41*, 886-896.

Kalish, M., Lewandowsky, S., & Kruschke, J. K. (2004). Population of linear experts: Knowledge partitioning and function learning. *Psychological Review*, *111*, 1072-1099.

Koh, K., & Meyer, D. E. (1991). Function learning: Induction of continuous stimulus-response relations. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *17*, 811-836.

Lewandowsky, S., Kalish, M., & Ngang, S. K. (2002). Simplified learning in complex situations: Knowledge partitioning in function learning. *Journal of Experimental Psychology: General*, *131*, 163-193.

McDaniel, M. A., & Busemeyer, J. R. (2005). The conceptual basis of function learning and extrapolation: Comparison of rule-based and associative-based models. *Psychonomic Bulletin & Review*, *12*, 24-42.

Navarro, D. J., & Perfors, A. (2009). Learning time-varying categories. In *Proceedings of the 31st annual conference of cognitive science society* (p. 414-424). Austin, TX: Cognitive science society.

Navarro, D. J., & Perfors, A. (2012). Anticipating changes: Adaptation and extrapolation in category learning. In N. Miyake, D. Peebles, & R. P. Cooper (Eds.), *Building bridges across cognitive sciences around the world: Proceedings of the 34th annual conference of the cognitive science society* (p. 809-814). Austin, TX: Cognitive Science Society.

Navarro, D. J., Perfors, A., & Vong, W. K. (2013). Learning time-varying categories. *Memory and Cognition*, *41*, 917-927.

Pearce, J. W. (2007). Psychopy - psychophysics software in python. *Journal of Neuroscience Methods*, *162*, 8-13.