

Lecture Notes

Lenard Dome

June 25, 2025

Contents

1	Why model?	2
2	Definitions of Computational Modeling	2
3	An example: the Rescorla-Wagner model	3
4	Discrepancy functions	5
5	Parameter estimation	6
6	Parameter recovery	8
7	Model comparisons	9
8	Model recovery	10
9	Toolbox	11
9.1	Model parameters	12
9.2	Model construction	13

1 Why model?

Let us look at a problem first. Imagine it is Friday night, you are done for the week, and I mean you are wasted. You cannot be bothered to cook, so you order a pizza.

2 Definitions of Computational Modeling

So, this is a crash-course slash half-day workshop on computational modeling, so we must begin with some definitions. This will be our first problem. As you might have observed, definitions of things become somewhat vague when they become popular. I broke this down into three distinct groups that I think are useful for us to consider:

1. Computational modeling is related to the formal expression of theories of cognition. A formal psychological model, as defined in this group, is a mathematically specified theory that is implemented as a computable algorithm for the purpose of simulating some aspect of human, or animal, behavior. This is the straining, where the formalisation helps us posit mechanisms of cognition, that we then implement to test whether our theories correspond to the real world. Two examples here are theories about the representations underlying categories (exemplars vs prototypes) and theories of attention and learning (e.g. attention is allocated to reduce prediction error).

2. The next category is more about measuring some aspects of behaviour. In this case, computational models are defined to measure some form of psychological construct (e.g. bias-free accuracy in recognition memory tasks). Two examples there are models based on Signal Detection Theory (SDT) to measure how well confidence judgements are aligned with the actual performance of the participant; and for example establish reference points via normative modeling, which is then used to measure individual differences in cognitive performance or any domain really.

3. The last category, where computational modeling as a phrase appears, is about the use of computational models as in statistical modeling. For example, computational models as statistical models help us cluster participant strategies into finite set of groups, analysing generalization gradients across continuously varying stimuli, or quantifying relationships between variables in a dataset.

One thing you might notice when you engage with the literature is that these three categories are overlapping. For example, a model that is used to measure some aspect of behavior can incorporate (sometimes implicitly) some theoretical assumption about cognition. Similarly, statistical models can be used to formalize theories or measure constructs.

We are largely concerned with the first category today, but our toolbox (the one we will be getting acquainted today) covers mainly the first two even though you can do all three.

3 An example: the Rescorla-Wagner model

So, the best way to continue is to give you an example of a computational model. The Rescorla-Wagner model is a classic model of associative learning, which is a type of learning that occurs when an organism learns contingencies between stimuli and outcomes. I pick two reasons for this example: first, the Rescorla-Wagner delta rule equation is probably the most influential equation in psychology, and second, at least three groups of people came up with the idea independent from each other (one from psychology animal cognition, one from artificial intelligence, and one from electrical engineering).

Now I will not go into a history much deeper, but let me just briefly set the context here. Models of learning at the time were largely single linear operators, such as Bush and Mosteller's model, where it was assumed that the strength of the association between a stimulus and an outcome is not just a linear function of the number of times the stimulus and outcome are presented together, but a decelerating function, where learning tends to slow as it proceeds. This is depicted in a famous learning curve you see there, where learning moves towards an upper limit, delta, as the number of trials increases. That is the prediction error, the model expressed by Bush and Mosteller, capturing the idea that the strength of an association between a stimulus and an outcome first increases rapidly, and then slows down with time (as prediction error decreases according to the model). This increase is also modulated via a learning rate parameter, which is a constant that determines how quickly the association is learned.

Now this is not a good model of how people or even animals learn. How do we know this? I give you one reason from Psychology, but you can probably find some in other fields as well.

Enter Blocking. A mouse first repeatedly observes that a light is followed

by food. It learns to salivate to the light. Next a tone is presented at the same time as the light, again followed by food. This also happens multiple times. Finally, the tone is presented alone, without the light.

What does the Bush-Mosteller system predict in this case?

Well, if you pair light with food often enough, the association will max out, and then stop. Now, when you present the tone and light together, the association to the light won't increase any further, but the association between tone and food will again, max out. So, when the tone is presented without the light, the mouse expects food and so salivates.

This is not what happens. In fact, little salivation occurs to the tone alone. This is how the phenomenon gets its name - the association of the tone to the food seems to be blocked by the pre-existing association of the light to the food.

Rescorla and Wagner (1972) proposed a new model to account for this. They made a very slight modification to the Bush-Mosteller model.

the Rescorla-Wagner theory takes the position that learning should stop when the sum of the input activations matches the teaching signal. In other words, when the outcome - food - is fully predicted. This idea allows it to predict blocking.

Let's see how that works. When the food and light are presented together, the association between them rises to one and stops, just like in Bush-Mosteller.

However, now the tone comes along. The food is already fully predicted by the light, and so from the perspective of Rescorla-Wagner, there's nothing to learn. The weight between tone and food remains at zero. This is often described as error-correcting learning - we only learn to the extent we need to in order to avoid errors.

I included a more general form here, which is what we usually mean nowadays when we say the delta rule.

We will use a version of this learning rule in our worksheets later, but I want to make the point that building models are hard. The simple modification they made was not mathematically discovered, like they did not stumble upon it by accident through rearranging equations, but was a result of implementing some hypotheses about how learning works. So we will do something similar here, apart from the fact that we will not invent a new model, but montage together a model from existing components. We will have a task, and we will have a model. So the remainder of this lecture, I will largely talk about what we do after the model is built.

4 Discrepancy functions

Now, we have a model, and we have some empirical findings. What is next? Back in the day, just like in the Rescorla-Wagner example, people plotted or displayed the results of model simulations and empirical data side by side, and then they looked at the graphs and said: "yep, looks good". This works very well if the data you have is robust and very good at distinguishing between models, and if the model is simple, but it is not a very good way to quantify how well the model approximates the data, especially if there are alternatives.

Here come discrepancy functions. Discrepancy functions are mathematical functions that quantify the difference between two sets of datapoints. They are often used to compare the predictions of a model with observed data, and they can help us determine how well a model **fits** the data.

So, we have a problem at our hand. We have the data, Y , the model M , and we want to know how close the two gets.

(next slide)

The first set of discrepancy functions are the ones that are based on the sum of squared differences. These functions are often used in regression analysis, where we want to minimize the difference between the predicted values and the observed values. The most common form is the least squares error, which is simply the sum of the squared differences between the predicted values and the observed values. You can also use them in computational modeling to capture how well the model is doing. These are non-parametric methods, which is great, because they do not make any assumptions about the distribution of the data. However, they are sensitive to outliers, which can skew the results.

But most often, we will use something else, called likelihoods. There are a few reasons for using them. First, models often assign probabilities to the data. On a given trial, in an experiment, there is a certain probability of observing a certain outcome. Like the probability of picking one from two alternatives (one red pill and one blue pill). The insight that powers likelihoods is that there is a function that can relate this model-assigned probability to the observed data.

For example, if the model predicts that the probability of picking the red pill is 0.7, and the observed data shows that the participant picked the red pill, then the likelihood of this outcome is 0.7. If the participant picked the blue pill, then the likelihood is 0.3 (which is a bad fit). The likelihood

function captures this relationship between the model's predictions and the observed data.

Of course, this is all conditional on the model, M , and its current parameters, θ . So, we can write the likelihood function as a function of the model and its parameters, which is what we do here.

Now there is a distinction between two ways these probability functions can be defined. The first, on the left, is when the model actually predicts a probability distribution over the data. This is the case for a Wald Drift Diffusion model, for example. The second, on the right, is when the model predicts a single value for each data point, such as the probability of a given outcome (0.7 chance of taking the red pill). We will then need to define a probability distribution over the data, for example a Bernoulli distribution for binary outcomes (like picking one of two pills).

Today, we will focus on the second type, where the model predicts a single value for each data point, and we will use a Bernoulli distribution to define the likelihood function. This is because it is easier to implement and understand, and it is sufficient for our purposes.

So, I outline here how this is done for a binary outcomes. We have our model assigned probabilities given the model parameters, θ , and the data, Y . Using the Bernoulli distribution, we can define the likelihood function as the product of the probabilities of the observed outcomes. This is done by taking the product of the probabilities for each data point.

5 Parameter estimation

So, we have a model, and we have a discrepancy function that quantifies how well the model fits the data. What is next? The next step is to estimate the parameters of the model that best fit the data. We are (for our current purposes) not interested in all possible model output values for all possible parameter values, but rather want to know the parameter values that make the model fit the data best. We are trying to find the peak of the likelihood function, the maximum of the likelihood, given the data and the model. This is called maximum likelihood estimation (MLE).

Most optimisation procedures rely on minimising an objective function, so we will actually minimise the negative log likelihood, which is a common practice and is shown here. There are all sorts of advantages, mainly numerical stability, but it is enough for you now to know that this is what we will

do.

So now we have the goal, what is the assumption here? Models that approximate the data well can tell you something about the utility of the model. Essentially, we assume that the model is a good approximation of the data, and that the parameters of the model can be estimated from the data. This is a strong assumption, but it is often necessary to make in order to use computational models in psychology.

Why do I say that? Because there are caveats to this. Goodness-of-fit does not capture a lot of Something I investigated throughoughly in my own research.

Let us start with flexibility. A standard rule of thumb says that the number of free parameters provide an upper bound for model flexibility... if a theory has five orthogonal free parameters, then it will be able to fit exactly any five data points; if the parameters are not orthogonal, however, the number of data points the theory can fit exactly is less. Now we will talk about whether the number of parameters are a good indicator of model complexity - briefly, it is not, and it has been empirically shown multiple times that the number of parameters is often not a good indicator of model flexibility (or complexity) when we talk about cognitive models.

Goodness-of fit also doesn't tell you about how diagnostic your experiment is. Simply put, it requires us to understand what are all the predictions the model can make in a given experiment. Are there things the model does not predict? Can the experimental design allow for those results as well? It also requires us to properly explore variability in the data, and also to show it. Now we usually do this by fitting on a subject-level (fit the model against each subjects data separately).

The last point I want to make here is that goodness-of-fit does not tell you about what is sometimes called as "plausible" falsifiability - are there any results the model could not fit? From a slightly different approach, something I prefer, if the model can account for the data as a function of it being able to account for all possible data, then it is a bad model. A model that is able to produce all logically possible results, observed and unobserved, can never fail to accommodate a result. Such an overly flexible model is inadequate through a lack of specificity.

6 Parameter recovery

After our brief walkthrough from model construction and parameterisation, parameter estimation, and discrepancy functions, I would like to talk a bit about a way to evaluate how reliable your model is before any data has been collected. This is called parameter recovery.

Mathematical models often reflect real-world systems and their parameters have biological, chemical, or physical interpretations, and not identifying these parameters can result in ambiguous interpretations. Parameter recovery is a method to assess the reliability of parameter estimates by simulating data from a model with known parameters and then attempting to recover those parameters from the simulated data.

This method comes in handy when you need to determine whether your model can accurately estimate the parameters of interest - you have to determine whether model parameters are identifiable, meaning that a given data can be only capture by unique parameter values.

Parameters are identifiable when there is only one possible value given the data.

Parameters are not identifiable, when there are multiple possible values that can explain the data. This can happen when the model is too flexible, or when the data is not informative enough to distinguish between different parameter values. As models become more complex, like some of the more complex recurrent networks and cluster-based categorisation models I dealt with, evaluating this becomes difficult, because the same dataset can be produced by multiple parameter sets, each corresponding to different underlying representations in the model - even though the model can produce the same data with different parameter sets, its internal state can differ, which corresponds to different predictions about what underlies a certain behavioural response. For the model you will use in the workshop, this is not a problem, because the model is simple enough to not have to worry about it.

This here shows the result of a prototypical parameter recovery procedure. We have two parameters, learning rate and choice stochasticity - don't worry too much about what they mean, you will be explained in the worksheets.

The Pearson correlation is quite high, which means that parameters are identifiable.

7 Model comparisons

So far, we only looked at ways that quantify how well the model does by itself. Now we shift towards looking at how well the model does compared to other models. Model comparison is crucial, especially because we often want to exclude potential theories in order to arrive at an adequate model of whatever we are trying to model.

The first issue I want to talk about is **overfitting**. You have probably heard this word many times and it is often used in the context of machine learning, but it is also relevant for computational modeling. Overfitting occurs when a model fits the training data really well but may perform poorly for independent data. Sometimes it is due to model complexity, but it largely depends on how you view complexity or how you measure it. It can also happen due to overtraining models, which is more of a problem in larger neural networks. But the problem is that the model does not generalize well to new data, and it is therefore not a good representation of the underlying system.

On the right here, this is demonstrated really well. You have three models with increasing complexity (increasing number of parameters) trying to predict life satisfaction scores as a function of years in marriage. You can see that the more complex models fit the training data better, but they do not generalize well to new data - no one expects satisfaction with life to drop to 0 after 10 years of marriage, but the most complex model predicts that. So, what does this mean for us? If we stick with the complexity dimension, we have a push towards simpler models, because they are less likely to overfit the data. But we also have a push towards less model parameters. Usually in the case of less than the number of data points.

This is the way we usually approach the problem of overfitting - through complexity as measured by the number of parameters. Unfortunately, the flexibility added by a free parameter depends on the details of the theory (cf. $\alpha x + \beta x$ with $\alpha x + \beta$; both have two parameters, but the latter is more flexible). The only accurate way to "allow" for the flexibility of a theory, as far as we know, is to determine what the theory predicts, and for that you have to simulate. Some methods try to address this shortcoming by looking at incorporating variance explained by each parameter into the penalty term. These include information criterion measures that generalize the goodness-of-fit measures across the model's parameter space, assuming that the data is large enough and the parameters are distributed according to Jeffrey's prior.

And other measures that try to add measures on how well each parameter

Another important aspect of model comparison is distinguishability. The problem we are facing here is that we often have multiple models and want to be able to identify the one that is better at explaining the data. Model recovery is a method to assess whether we can identify which model generated the data. This is done by simulating data from each model and then fitting the models to the simulated data. Then we can count how many times each model came out on top when fitted against the simulated data, which gives an indication of how well the models can be distinguished from each other. We will talk about it in a second in more details.

Some extra techniques I want to mention because I think it is important to be aware of them, even though we will not cover them, are landscaping and parameter space partitioning. Landscaping is a way to extend model recovery by comparing model fits against each other visually. It is a way to check how well one model is preferred over the other in a given comparison. Parameter space partitioning is a little bit more sophisticated. In this approach, we turn model predictions into a finite set of countable discrete outcomes, and we use those to identify regions in the parameter space that correspond to different discrete outcomes (such as different response patterns in a categorization task). This method gives an interesting measure of complexity, that is based on actual model outcomes as opposed to the number of parameters.

8 Model recovery

We will talk about model recovery in more detail, because this is good and standard practice that most people should do before data collection or when planning to develop a research programme or a series of studies that involve model comparison.

In a model recovery, we extend what we learnt about parameter identifiability to models. In this scenario, we generate data from a model with known parameters, and then we fit multiple models to the data. We repeat this process for all data generating model many times, and estimate how likely that the model wins over the alternatives for a given dataset. So, we acquire what is often called confusion matrix, see on the right. Each cell gives you a probability of a model given the data, and the rows are the data generating models. The diagonal gives you the probability of the true data-generating model winning over the alternatives, and the off-diagonal cells give you the

probability of the model is losing to the alternatives.

If the models are distinguishable, then the diagonal cells will be high, and the off-diagonal cells will be low. If the models are not distinguishable, then the diagonal cells will be low, and the off-diagonal cells will be high or comparable. Now it is often the case that models are misspecified and become indistinguishable, like when they share the same parameters or parts of their architecture. Or it could also be the case that the experimental design is not complex enough to distinguish between the models. In this case, we have to return to the drawing board and either change the model or the experimental design.

Now beware that what we mean by confusion matrix here is not what we mean by confusion matrix in machine learning. It is an important distinction, but they can be related under certain conditions.

9 Toolbox

So, we have now covered the basics of computational modeling, and we have a good understanding of how to build models, estimate their parameters, and compare them. Now we will turn our attention to the toolbox that we will use in the workshop.

I organised these sections to correspond with the exercises that we will be doing and given that we have now covered most of the theoretical basics, we can move on to relate the theory to practice.

First off, the toolbox is new, fresh out of the box. Which means that you will have an impact on what it will become. Thank you for that. Also, the toolbox is not a monolithic entity, but rather a collection of modular and interoperable tools that can be used for different purposes. Some of the tools are more general, while others are more specific to certain types of models or data. The toolbox is designed to be flexible and extensible, so you can add or create your own tools or modify existing ones to suit your needs.

The toolbox also relies on pandas, numpy, and scipy, so if you have encountered those packages, or if you want to do data analysis in python, you can easily generalise what you already know to what we will be doing. And generalise what you learn here to some of the other packages.

One reason for doing this is to avoid reinventing the wheel. If there is a good implementation of gradient search, why try to write your own?

But let's not spend too much time here. Let's move on to what we are

doing during the exercises. The idea here is to give you a brief overview of what might be a simple computational workflow with some optional extras that extend it.

Start with the toolbox. As I said before it is a collection of tools that are aimed to assist you in different parts of modeling. It is composed of modules (standard in a python library) that you can call from the standard cpm library to assist you in different modeling steps. The very dense figure you see here depicts how these might work. Each rectangle is a different module, and the arrows show how they are connected to each other. Each group of module indicated by the coloured dashed lines help you with a distinct part of the modeling process. The big arrow on the left shows you a simple modeling workflow, where you start with a model, then you estimate its parameters, draw some inferences on the group-level with techniques we discuss later, and finally you simulate with the estimated parameters. Model and parameter recovery lies outside of this arrow, but similarly involves interaction between modules that are not shown here, but we still do.

So, we have a set of modules that help you build your model, a set of tools that help you estimate the parameters of your model, and a set of tools that help you do hierarchical modeling. We also have a set of tools that help you simulate data from your model, which is useful for model recovery and parameter recovery.

9.1 Model parameters

Previously, we started with building a model. We start our toolbox-walkthrough there as well. The first thing we do here is parameterisation of the model. What are the parameters of the model? These parameters will tune the cognitive processes we formally express as equations here.

The toolbox provides two robust and important modules for this. The first is the `cpm.generators.Value` module, which provides a set of tools to help define each parameter. The second is the `cpm.generators.Parameters` module, which takes each individual parameters and other variables you need for your model and stores it in one convenient python object. What does this mean? It means that you can define your model parameters in a structured way, that allows you to easily change the parameters of your model without having to rewrite the code. In addition, it gives you several useful methods to manipulate the parameters, such as adding, removing, or modifying them, but more importantly, helps you sample parameters from a prior distribution,

calculate posterior densities, update priors, and archive the parameterisation of your model for things such as lab reports.

9.2 Model construction

The next step is to build the model. To that effect, the toolbox includes a library of algorithmic components, that you can use to montage your model together.