Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Model selection, comparison, and validation

Kate Nussenbaum | Vasilisa Skvortsova | Johanna Habicht

flux

Modeling Flux | Part 4 | Model Comparison & Recoverability

Developmental
Computational
Psychiatry lab

HARTLEY LAB

**1. Develop task and model(s)**

Task design ⟷ Model design

**2. Fit model(s) to data**

Model selection → Parameter estimation

**3. Validate results**

Model recovery and comparison | Parameter recovery | Posterior predictive checks

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

**3. Validate results**

| Model recovery and comparison | Parameter recovery | Posterior predictive checks |

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

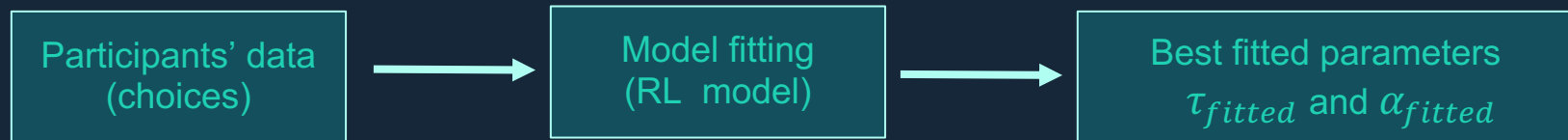What do all these procedures have in common?

**MODEL SIMULATION**

3. Validate results

Model recovery and comparison

Parameter recovery

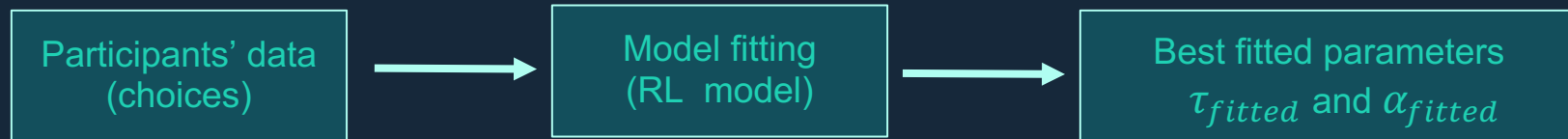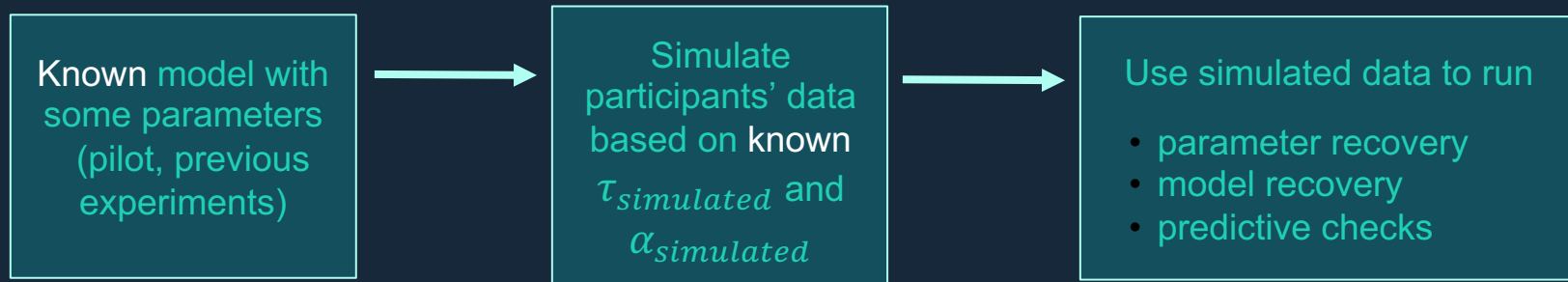Posterior predictive checks

flux

# Model simulation

| Participants' data (choices) | → | Model fitting (RL model) | → | Best fitted parameters $\tau_{fitted}$ and $\alpha_{fitted}$ |
|---|---|---|---|---|

→ **You start with the data and unknown model & parameters**

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

# Model simulation

Developmental
Computational
Psychiatry lab

HARTLEY LAB

| Participants' data (choices) | → | Model fitting (RL model) | → | Best fitted parameters $\tau_{fitted}$ and $\alpha_{fitted}$ |
|---|---|---|---|---|

→ **You start with the data and unknown model & parameters**

| Known model with some parameters (pilot, previous experiments) | → | Simulate participants' data based on known $\tau_{simulated}$ and $\alpha_{simulated}$ | → | Use simulated data to run<br>• parameter recovery<br>• model recovery<br>• predictive checks |
|---|---|---|---|---|

→ **You start with the known model & parameters to generate "fake" data**

flux

# Parameter recovery

*How reliable are model parameters?*
*How do parameters change relative to one another?*

→ **We need to perform parameter recovery checks.**
→ **For example, with current task and 1 LR model:**

- Recover softmax decision temperature
- Recover learning rates

HARTLEY LAB

Developmental
Computational
Psychiatry lab

flux

# Parameter recovery

**Steps:**
1. **Fit the model to behavior and define parameters' range (average/median/min-max)**
2. **Simulate the model varying one of the parameter values while keeping other parameters fixed**
3. **Fit simulated data with the same model used for the simulation**
4. **Compare true and recovered parameters.**

# Parameter recovery

**Example:**

**2-arm bandit task with binary outcomes (reward, no reward)**
**Model with 1 learning rate and softmax decision rule**

N = 30 subjects

T = 100 trials, $p_{reward} = 0.8$

Best fitted model group parameters:

chosen learning rate $\alpha_{chosen} = 0.5$,

softmax temperature $\tau = 0.3$

flux

# Parameter recovery

For softmax temperature $\tau$

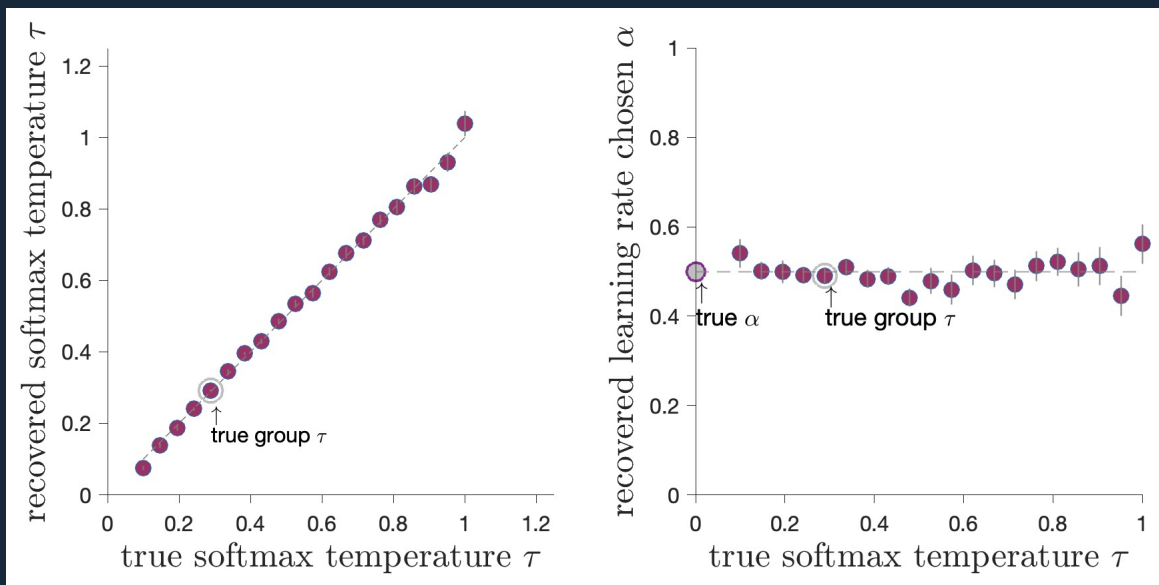Simulate $\tau \in [0.1, 1]$ and $\alpha_{chosen} = 0.5$

# Parameter recovery

For softmax temperature $\tau$

Simulate $\tau \in [0.1, 1]$ and $\alpha_{chosen} = 0.5$

# Parameter recovery

For softmax temperature $\tau$

Simulate $\tau \in [0.1, 1]$ and $\alpha_{chosen} = 0.5$

→ **Recovery shows problem with high values of the softmax inverse temperature:** $\uparrow \tau$ →**more exploration** → **less stable estimates**

→ **True values for both parameters might be overestimated**

→ **Try: increase the number of trials:** $100 \rightarrow 300$

HARTLEY LAB

Developmental
Computational
Psychiatry lab

flux

# Parameter recovery

For softmax temperature $\tau$

Simulate $\tau \in [0.1, 1]$ and $\alpha_{chosen} = 0.5$

# What is "good" recoverability?

→ No hard rule or accepted values

→ Experimental design and nature of the data: example laboratory well-controlled study vs. online experiment with much noisier data

→ Simulations help setting a benchmark & optimise the design

      → Vary number of trials
      → Vary parameter ranges

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

# Model recovery

→ Questions about cognition (and developmental change in cognitive processes) can be addressed through examining **parameter estimates** from a single model

flux

# Model recovery

→ Questions about cognition (and developmental change in cognitive processes) can be addressed through examining **parameter estimates** from a single model and/or by **comparing *different* models.**

flux

# Model recovery

→ Questions about cognition (and developmental change in cognitive processes) can be addressed through examining **parameter estimates** from a single model and/or by **comparing *different* models.**

→ In the domain of reinforcement learning, different models typically formalize different **value-updating processes** or **choice functions**.

flux

# What if there are multiple plausible models of behavior?

→ **Typically, more than one hypothesis about behavior can be formalized algorithmically**

# What if there are multiple plausible models of behavior?

→ **Typically, more than one hypothesis about behavior can be formalized algorithmically**

→ **For example, with task described earlier:**

1. One learning-rate model

2. Decay model

3. Null model

# Defining different models

1. One learning-rate model – 2 parameters $(\tau, \alpha)$

   Single learning rate scales prediction errors

# Defining different models

1. One learning-rate model – 2 parameters $(\tau, \alpha)$
   ## Single learning rate scales prediction errors

2. Decay model – 3 parameters $(\tau, alpha_{initial}, \eta)$
   ## Learning rate decays over time

$$\alpha_{decay} \; at \; t = \frac{\alpha_{initial}}{1 + \eta * trial_{t-1}}$$

# Defining different models

1. One learning-rate model – 2 parameters $(\tau, \alpha)$
   
   Single learning rate scales prediction errors

2. Decay model – 3 parameters $(\tau, alpha_{initial}, \eta)$
   
   Learning rate decays over time

3. Null model – 0 parameters
   
   No learning — random choice on every trial

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

# Model comparison: determining which model best captures data

**Steps:**

1. **Fit all possible models to behavior**
2. **Compare indices of "fit"**

# Model comparison: determining which model best captures data

**Steps:**

1. **Fit all possible models to behavior**

2. **Compare indices of "fit"**

   → Take into account likelihood — *the probability of the observed choices given the algorithm* — AND penalize more complex models

# Model comparison: determining which model best captures data

**Steps:**

1. **Fit all possible models to behavior**

2. **Compare indices of "fit"**

   → Take into account likelihood — *the probability of the observed choices given the algorithm* — AND penalize more complex models

   → Two common metrics: AIC and BIC

# Model comparison: determining which model best captures data

**Steps:**

1. **Fit all possible models to behavior**

2. **Compare indices of "fit"**

     → Take into account likelihood — *the probability of the observed choices given the algorithm* — AND penalize more complex models

     → Two common metrics: AIC and BIC

          AIC: $2k - 2\ln(L)$

          BIC: $k\ln(n) - 2\ln(L)$

# AIC and BIC

AIC: $2k - 2\ln(L)$

BIC: $k\ln(n) - 2\ln(L)$

*k: number of parameters*

*L: max likelihood*

*n: number of observations*

→ **Smaller values are better**

# Finding the best-fitting model

# Finding the best-fitting model

## Likelihood term

# Finding the best-fitting model

# Finding the best-fitting model

# Are our models 'recoverable'?

→ **Critical to ensure that different models are actually distinguishable from one another, given the task design.**

# Are our models 'recoverable'?

→ Critical to ensure that different models are actually distinguishable from one another, given the task design.

→ Extreme example:

Imagine a task that involves 3 trials.

Can quantitatively assess model fit, but it's unlikely you will really be able to learn anything about the cognitive processes behind a participant's choices.

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How do we know whether our model-fitting results reflect reality?

→ *Problem:* No way to know the 'true' algorithm a participant used to make choices.

→ *Solution:* Simulate fake participants so that we *know* the algorithm that generated the choice data.
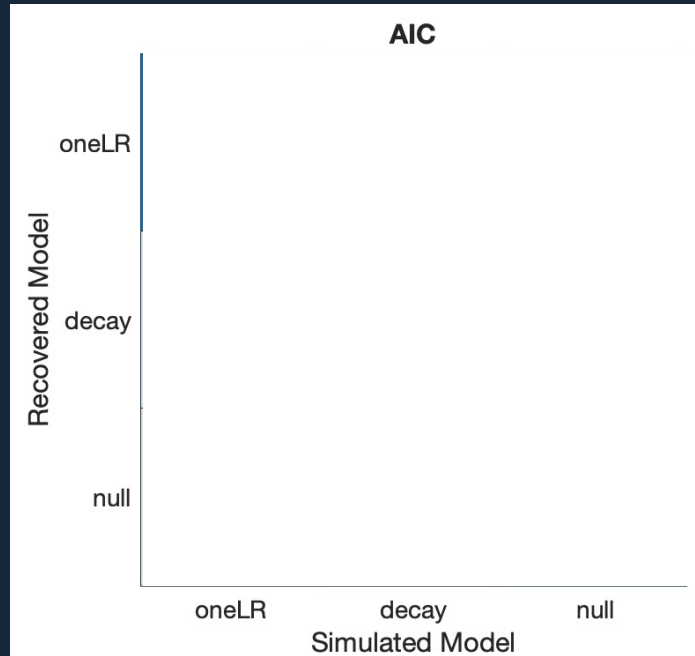
flux

# Model recoverability analyses

**Steps:**

1. Simulate data from all models.

2. Fit models to all simulated datasets.

3. Determine which model best fits each data set.

4. Determine the proportion of datasets for which the 'recovered' model matches the 'ground truth' model.
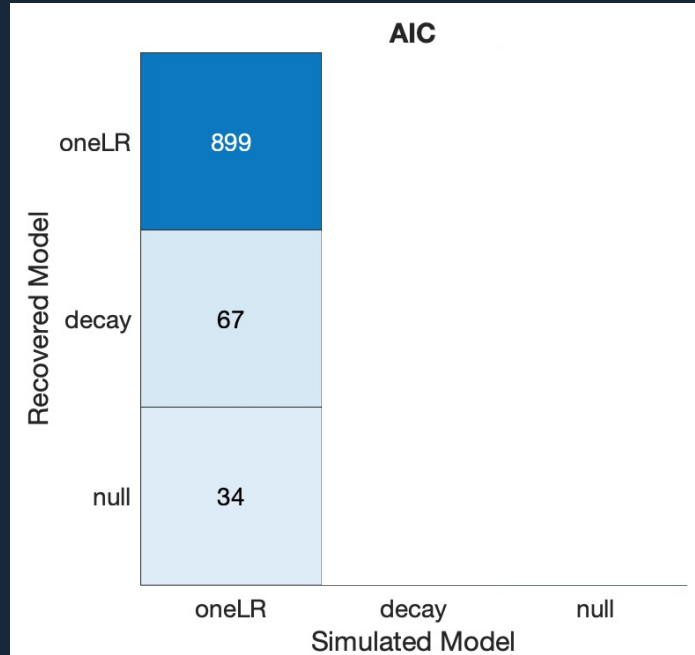
Developmental
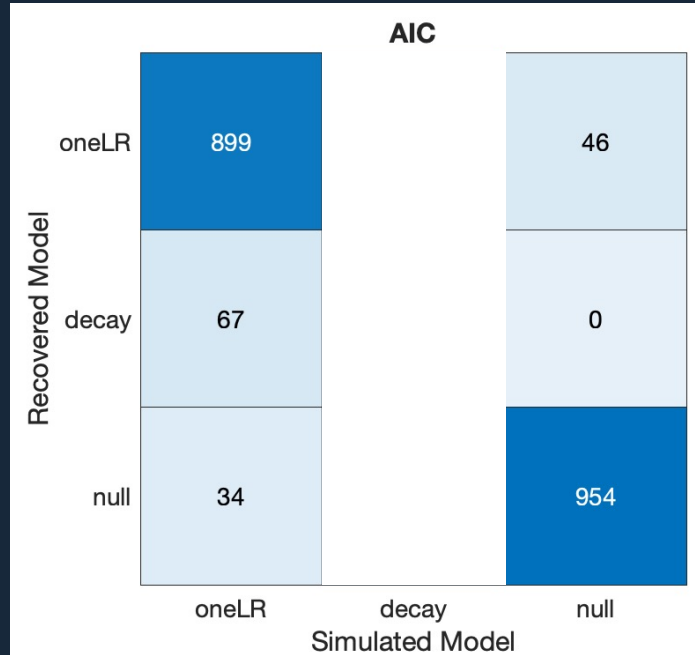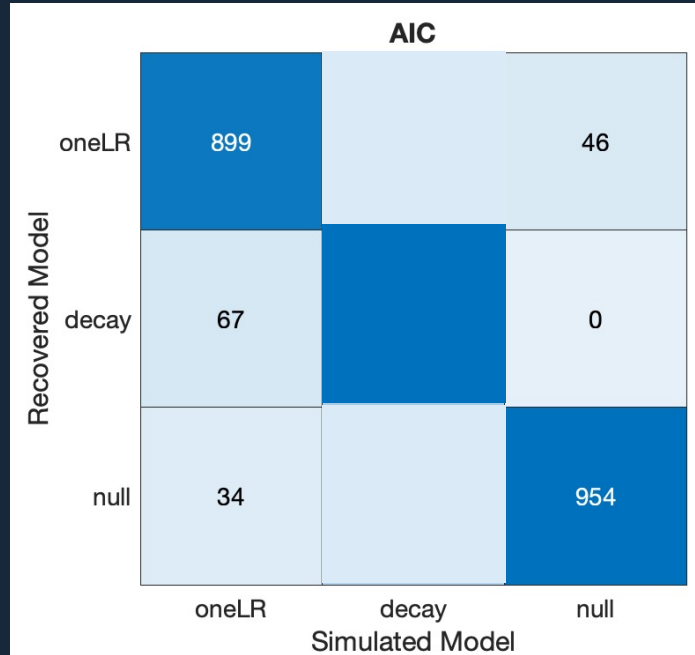Computational
Psychiatry lab

HARTLEY LAB

flux

# Model recoverability analyses: Confusion matrices

# Model recoverability analyses: Confusion matrices

# Model recoverability analyses: Confusion matrices

# Model recoverability analyses: Confusion matrices

# Model recoverability analyses: Confusion matrices

# Task optimization

*What if model recoverability is poor?*

→ **Change aspects of task design to improve it.**

flux

# Task optimization

*What if model recoverability is poor?*

→ **Change aspects of task design to improve it.**

→ Examples: Number of trials, number of stimuli, changes in reward probabilities, etc.

# Task optimization

*What if model recoverability is poor?*

→ **Change aspects of task design to improve it.**

→ Examples: Number of trials, number of stimuli, changes in reward probabilities, etc.

→ **Repeat.**

# Comparing task versions

## 20 trials



## 200 trials

# Comparing task versions

## 20 trials



**AIC**

|  | oneLR | decay | null |
|---|---|---|---|
| **oneLR** | 631 | 359 | 48 |
| **decay** | 20 | 83 | 2 |
| **null** | 349 | 558 | 950 |

Recovered Model / Simulated Model

## 200 trials



**AIC**

|  | oneLR | decay | null |
|---|---|---|---|
| **oneLR** | 904 | 496 | 59 |
| **decay** | 91 | 404 | 2 |
| **null** | 5 | 100 | 939 |

Recovered Model / Simulated Model

# Predictive performance and model checks

1.  **All models could be wrong — *the model comparisons are relative***

    ➔ Better "fit" does not guarantee the model reproduces behaviour

    ➔ Better "fit" does not guarantee the model could be recovered

# Predictive performance and model checks

1. **All models could be wrong — *the model comparisons are relative***

    → Better "fit" does not guarantee the model reproduces behaviour

    → Better "fit" does not guarantee the model could be recovered

2. **Generative performance of the model: how well the model can reproduce behaviour**

# Predictive performance and model checks

1. **All models could be wrong —** *the model comparisons are* **relative**

    ➔ Better "fit" does not guarantee the model reproduces behaviour

    ➔ Better "fit" does not guarantee the model could be recovered

2. **Generative performance of the model: how well the model can reproduce behaviour**

    ➔ Check model performance against behavioural data **qualitatively**

    ➔ Try to find a behavioural pattern that **dissociates** between the models

*The importance of model falsification (Palminteri et al., Trends Cog Sci 2017)*

# Generative performance of the model

**Steps:**

1. **Simulate the models with the best fitted parameters**
2. **Define a behavioural marker:  where models' behaviour won't generate the same predictions**
3. **Compare model performance to subjects' actual behaviour**

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Generative performance of the model

**Example:**

**Model 1:** *fixed* **learning rate and softmax decision rule**
**Model 2:** *decaying* **learning rate and softmax decision rule:**
the agent progressively decreases the update of the option values
and "ignores" the irrelevant non-rewarding events

**The agent is** *less likely* **to switch choice after a negative prediction error**

**Benchmark behaviour:** $P(switch)$ **after a** *negative* **prediction error**
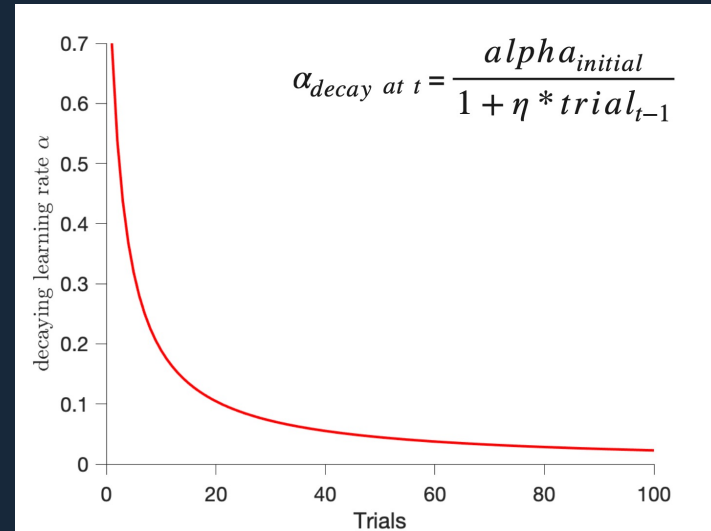
flux

# Generative performance of the model

**Example:**

N = 30 subjects who played 2-arm bandit task, T = 100 trials, $p_{reward} = 0.8$

Best fitted model group parameters:
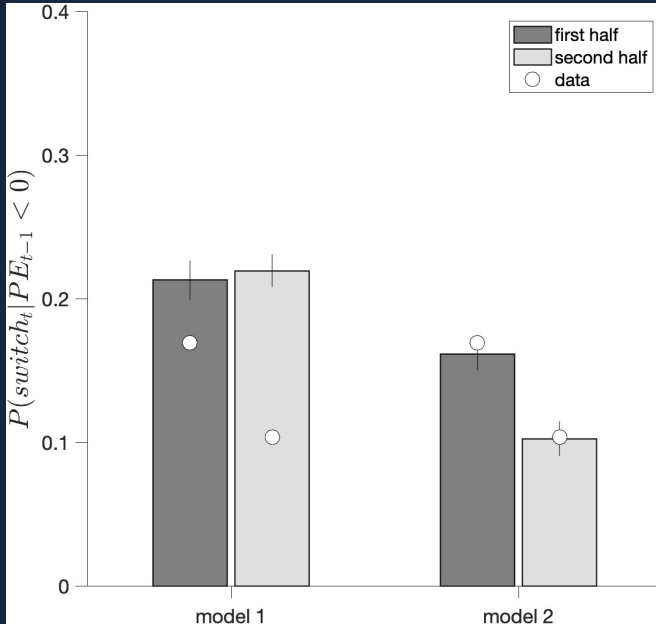
Model 1: chosen learning rate $\alpha_{chosen} = 0.7, \tau = 0.2$

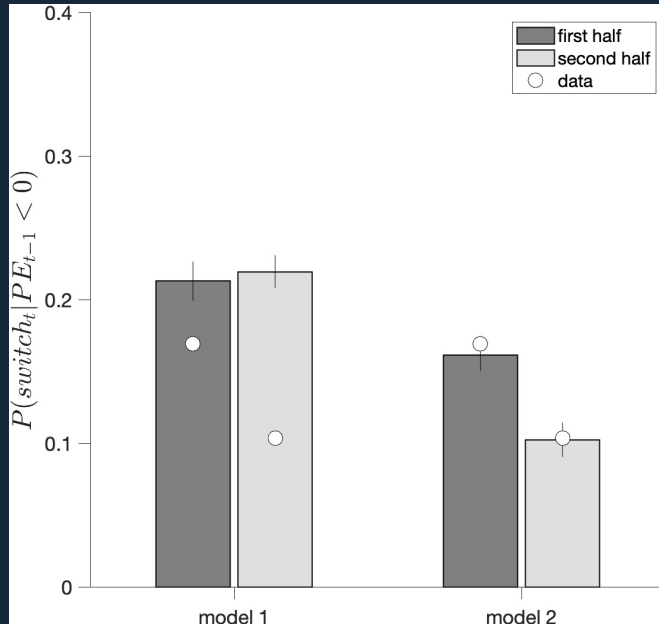Model 2: initial chosen learning rate $\alpha_{chosen} = 0.7$, decay parameter $\eta = 0.02$



$$\alpha_{decay\ at\ t} = \frac{alpha_{initial}}{1 + \eta * trial_{t-1}}$$

# Generative performance of the model

**Example:**

# Generative performance of the model

**Example:**



Only the model with *decaying* learning rate was capable of generating this behavioural pattern → the two models can be dissociated based on this behavioural marker.

# Model validation summary

→ It's *extremely important* to ensure parameters and models of interest are recoverable *prior* to data collection.

→ Simulation is a valuable tool to test modeling approaches.

→ Tasks and models should be developed and refined *concurrently.*

Q & A