Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How to Develop a Computational Model?
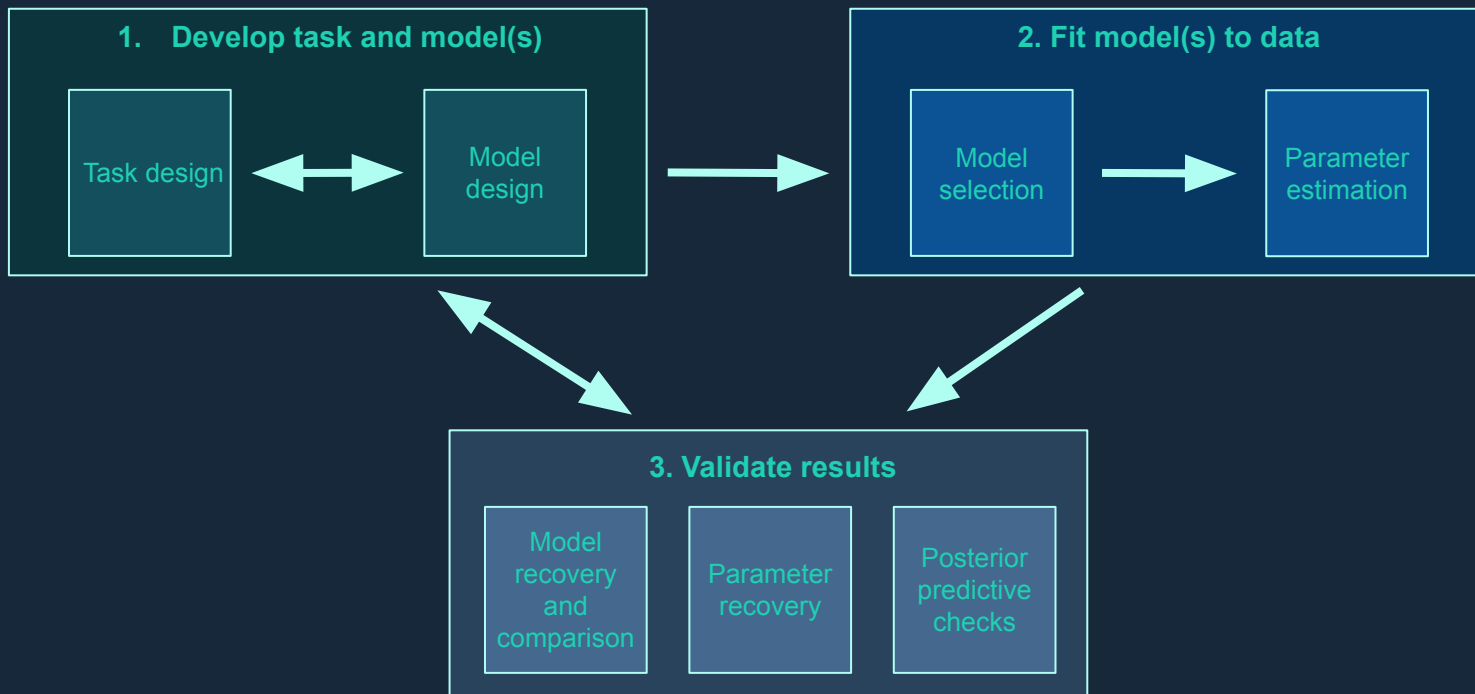
*"All models are wrong, but some are useful"*
*George E. P. Box*

Tricia Seow | Samuel Hewitt | Noam Goldway
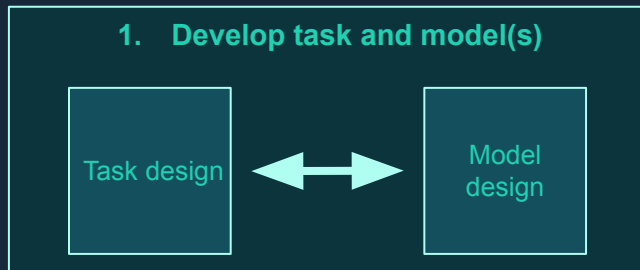
flux

# What we will cover:

➔ An example for how to select the **proper model with respect to a specific task design**

➔ The **Rescorla Wagner** model

➔ The concept of **learning rate**

➔ The concept of **temperature**

➔ What is a **"softmax"** function
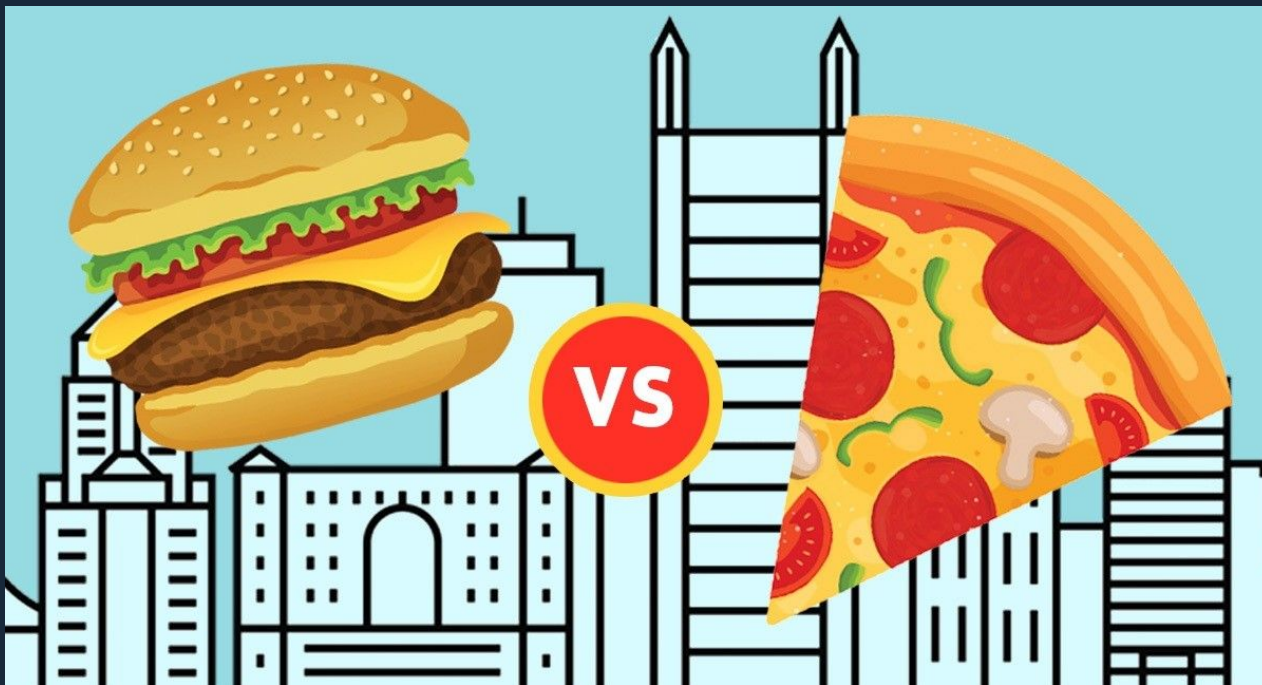
Tricia Seow | Samuel Hewitt | Noam Goldway

# Developing a computational model



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

Task design ↔ Model design

Tricia Seow | Samuel Hewitt | Noam Goldway

Tricia Seow | Samuel Hewitt | Noam Goldway

# Experimental task -

How do you maximise reward if you do not know which slot machine is better?
➔ Learn expected value of each slot machine
➔ Make the next choice based on values learnt

| Trial | Choice | Outcome |
|-------|--------|---------|
| 1 | **Right** | 0 |
| 2 | **Left** new choice | +1 |
| 3 | **Left** | +1 |
| 4 | **Right** | 0 |
| 5 | **Left** past experience | +1 |

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

# Modelling behaviour with RL

Value function: Rescorla Wagner model

$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$

Developmental
Computational
Psychiatry lab

HARTLEY LAB

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Prediction Error



$$V_{prepule} > V_{orange}$$

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Prediction Error



$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Prediction Error

$$V_t = V_{t-1} + \alpha(R_t - 0.5)$$

Tricia Seow | Samuel Hewitt | Noam Goldway

# Prediction Error

$$V_t = V_{t-1} + \alpha(1-0.5)$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Prediction Error



$$0.5$$
$$V_t = V_{t-1} + \alpha(1-0.5)$$

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

# Modelling behaviour with RL

## Value function: Rescorla Wagner model

$$V_t = V_{t-1} + \alpha (R_t - V_{t-1})$$

Value
(of the slot machine)

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

## Value function: Rescorla Wagner model

$$V_t = V_{t-1} + \alpha (R_t - V_{t-1})$$

Value
(of the slot machine) = Value on previous trial

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

## Value function: Rescorla Wagner model

$$V_t \quad = \quad V_{t-1} \; + \quad \alpha \; (R_t - V_{t-1})$$

| Value (of the slot machine) | = | Value on previous trial | ( Reward − Value on previous trial ) |

Prediction error
what you received - what you expected

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

## Value function: Rescorla Wagner model

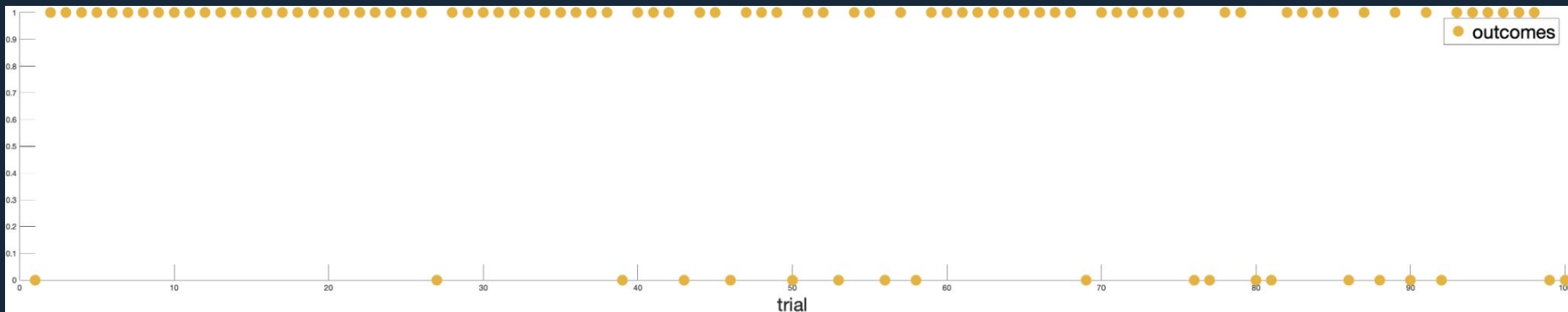$$V_t = V_{t-1} + \alpha (R_t - V_{t-1})$$

Value
(of the slot machine) **=** Value on
previous trial **+** Learning
rate **(** Reward **-** Value on
previous trial **)**

Prediction error
what you received - what you expected

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

Prediction error
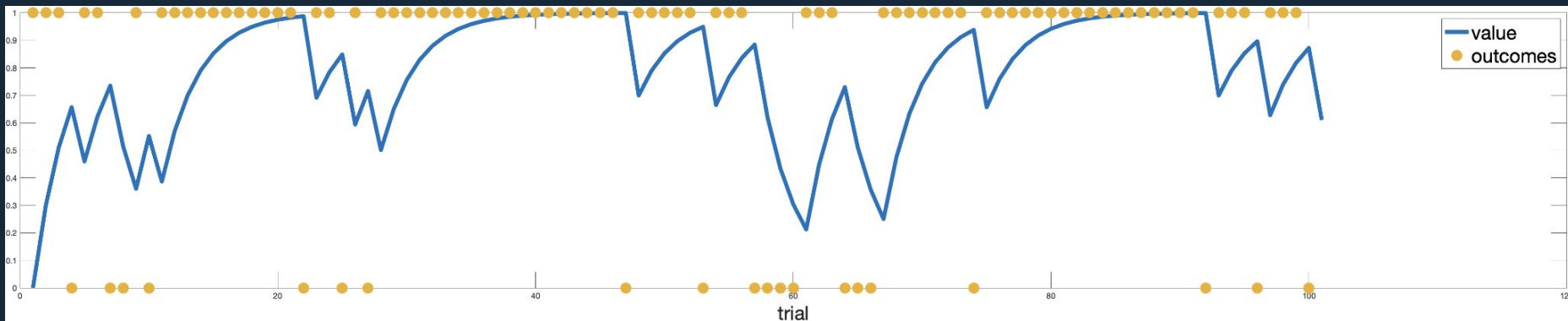what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \; ( \; \text{Reward} \; - \; \text{Value on previous trial} \; )$$



Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

Prediction error
what you received - what you expected

Value
(of the slot machine) = Value on previous trial + Learning rate ( Reward - Value on previous trial )



Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left( \text{Reward} - \text{Value on previous trial} \right)$$

| Trial 1 | ? | | | | |
|---------|---|---|---|---|---|

?

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

Value
(of the slot machine) = Value on
previous trial + Learning
rate ( Reward - Value on
previous trial )

| Trial 1 | ? | *initiation:* 0.5 | | | *initiation:* 0.5 |
|---------|---|-------------------|---|---|--------------------|

?

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

Value
(of the slot machine) **=** Value on
previous trial **+** Learning
rate ( Reward - Value on
previous trial )

| Trial 1 | ? | 0.5 | | 1 | 0.5 |
|---------|---|-----|---|---|-----|

Tricia Seow | Samuel Hewitt | Noam Goldway

# Modelling behaviour with RL

Prediction error
what you received - what you expected

Value
(of the slot machine) = Value on
previous trial + Learning
rate ( Reward - Value on
previous trial )

0.5

| | | | | | |
|---|---|---|---|---|---|
| *Trial* 1 | ? | 0.5 | + | *1* | *0.5* |

Tricia Seow | Samuel Hewitt | Noam Goldway

# Modelling behaviour with RL

Prediction error
what you received - what you expected

Value
(of the slot machine)  =  Value on previous trial  +  Learning rate  ( Reward -  Value on previous trial )

| Trial 1 | *1* | 0.5 | + | 0.5 |
|---------|-----|-----|---|-----|

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

Prediction error
what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left( \text{Reward} - \text{Value on previous trial} \right)$$

| | Value (of the slot machine) | Value on previous trial | + | Learning rate ( Reward - Value on previous trial ) |
|---|---|---|---|---|
| *Trial 1* | 1 | 0.5 | + | 0.5 |
| *Trial 2* | | 1 | + | 1 |

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

| | Value<br>(of the slot machine) = | Value on<br>previous trial | + | Learning<br>rate ( Reward - | Value on<br>previous trial ) |
|---|---|---|---|---|---|
| *Trial*<br>1 | *1* | 0.5 | + | *0.5* | |
| *Trial*<br>2 | | 1 | + | 1 | 1 |



1$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left( \text{Reward} - \text{Value on previous trial} \right)$$

|  | Value (of the slot machine) | = | Value on previous trial | + | Learning rate ( Reward - Value on previous trial ) | |
|---|---|---|---|---|---|---|
| *Trial 1* | 1 | | 0.5 | + | 0.5 | |
| *Trial 2* | | | 1 | + | (1 | - 1) |

Tricia Seow | Samuel Hewitt | Noam Goldway

# Modelling behaviour with RL

Prediction error
what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left( \text{Reward} - \text{Value on previous trial} \right)$$

| | Value (of the slot machine) | Value on previous trial | + | Learning rate ( Reward - Value on previous trial ) |
|---|---|---|---|---|
| *Trial 1* | 1 | 0.5 | + | 0.5 |
| *Trial 2* | | 1 | + | 0 |

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental Computational Psychiatry lab

HARTLEY LAB

flux

# Modelling behaviour with RL

$$\underset{\text{(of the slot machine)}}{\text{Value}} = \underset{\text{previous trial}}{\text{Value on}} + \underset{\text{rate}}{\text{Learning}} \left( \text{Reward} - \underset{\text{previous trial}}{\text{Value on}} \right)$$

Prediction error

what you received - what you expected

| | Value (of the slot machine) = | Value on previous trial | + | ( Reward - Value on previous trial ) |
|---|---|---|---|---|
| *Trial* 1 | *1* | 0.5 | + | 0.5 |
| *Trial* 2 | *1* | 1 | + | 0 |

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

Prediction error

what you received - what you expected

Value
(of the slot machine) = Value on previous trial + Learning rate ( Reward - Value on previous trial )

| | Value (of the slot machine) | Value on previous trial | + | Reward − Value on previous trial |
|---|---|---|---|---|
| *Trial* 1 | 1 | 0.5 | + | 0.5 |
| *Trial* 2 | 1 | 1 | + | 0 |

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental Computational Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

Prediction error
what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left( \text{Reward} - \text{Value on previous trial} \right)$$

| | Value (of the slot machine) | Value on previous trial | + | Reward - Value on previous trial |
|---|---|---|---|---|
| *Trial* 1 | 1 | 0.5 | + | 0.5 |
| *Trial* 2 | 1 | 1 | + | 0 |
| *Trial* 3 | | 1 | + | 1 |

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

| Value (of the slot machine) | = | Value on previous trial | + | Learning rate ( Reward - Value on previous trial ) |
| --- | --- | --- | --- | --- |

| | Value (of the slot machine) | Value on previous trial | | Learning rate ( Reward - Value on previous trial ) | |
| --- | --- | --- | --- | --- | --- |
| *Trial* 1 | *1* | 0.5 | + | 0.5 | |
| *Trial* 2 | *1* | 1 | + | 0 | |
| *Trial* 3 | | 1 | + | 0 | 1 |

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Modelling behaviour with RL

Prediction error
what you received - what you expected

| Value (of the slot machine) | = | Value on previous trial | + | Learning rate ( Reward - Value on previous trial ) |

| | Value (of the slot machine) | Value on previous trial | + | Value on previous trial |
|---|---|---|---|---|
| Trial 1 | 1 | 0.5 | + | 0.5 |
| Trial 2 | 1 | 1 | + | 0 |
| Trial 3 | | 1 | + | -1 |

Tricia Seow | Samuel Hewitt | Noam Goldway

# Modelling behaviour with RL

Prediction error

what you received - what you expected

$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left( \text{Reward} - \text{Value on previous trial} \right)$$

| | Value (of the slot machine) | Value on previous trial | | Value on previous trial |
|---|---|---|---|---|
| *Trial* 1 | 1 | 0.5 | + | 0.5 |
| *Trial* 2 | 1 | 1 | + | 0 |
| *Trial* 3 | 0 | 1 | + | -1 |

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental Computational Psychiatry lab

HARTLEY LAB

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Modelling behaviour with RL

Prediction error
what you received - what you expected

Value
(of the slot machine) = Value on
previous trial + Learning
rate ( Reward - Value on
previous trial )

$$V_t = V_{t-1} + \alpha (R_t - V_{t-1})$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How much should we learn?

## What happens if we manipulate learning rate?

$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$



Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How much should we learn?

## What happens if we manipulate learning rate?

$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$



α=0.1



α=0.3

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# Is low learning rate always better?

$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$

→ Depend on the statistics of the environment
- Low volatility-> low $\alpha$ is better
- High volatility-> high $\alpha$ is better

Tricia Seow | Samuel Hewitt | Noam Goldway

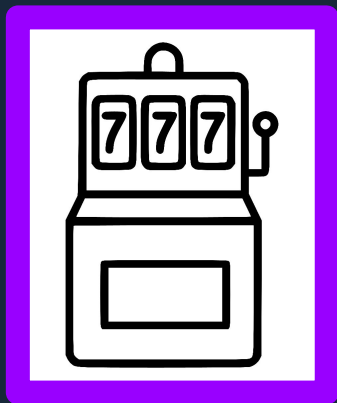https://github.com/DevComPsy/2021FluxCompModellingWorkshop

flux

# What did we learn so far

➔ What are multi arm bandit tasks
➔ How RL and, specifically Rescorla Wagner model can help us to 'solve' such problems
➔ Expected value
➔ Prediction error
➔ High vs low learning rate

Tricia Seow | Samuel Hewitt | Noam Goldway

https://github.com/DevComPsy/2021FluxCompModellingWorkshop

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?



*80% reward*

*20% reward*

← learnt via trial and error →
(value function)

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?



*80% reward*

Maximise rewards

→ Pick slot machine with largest likelihood of reward

→ Exploit

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?



Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

Choose **purple** is better

Choose **orange** is better
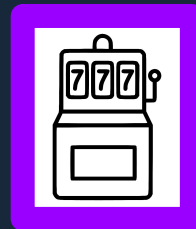


Tricia Seow | Samuel Hewitt | Noam Goldway

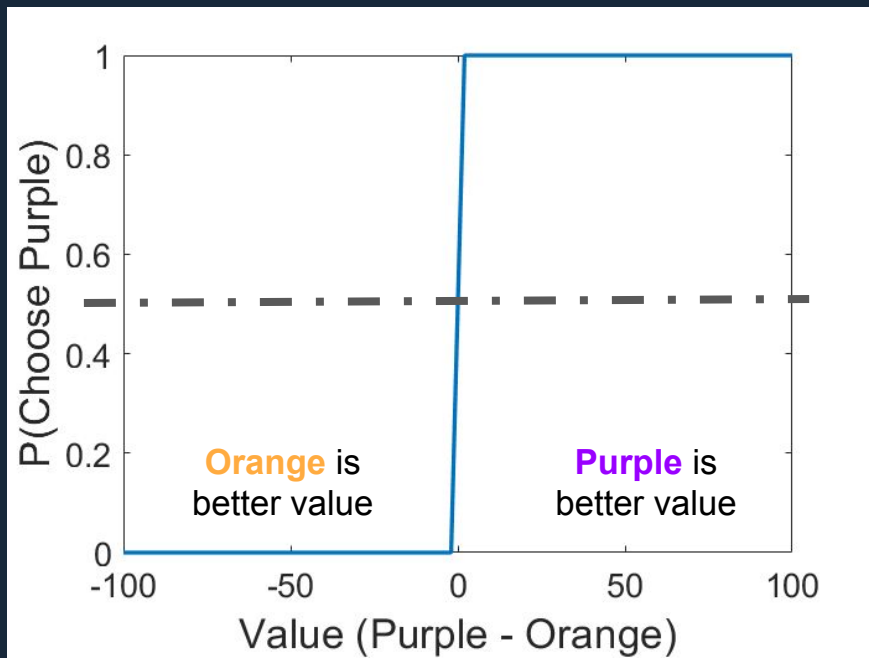# How should we choose?

Choose **purple** is better

Choose **orange** is better



*Exploit*

➔ Choose slot machine when reward is better than the other

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

## Try other options

➔ Sample the outcomes of the other slot machine

➔ Explore



*20% reward*

Tricia Seow | Samuel Hewitt | Noam Goldway

# How should we choose?

Choose **purple** is better

Choose **orange** is better



*Explore*

→ Choose slot machine equally

Developmental
Computational
Psychiatry lab

HARTLEY LAB

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

Tricia Seow | Samuel Hewitt | Noam Goldway

# How should we choose?

Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

## Low temperature
➔    Choices are less noisy



Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

## Low temperature

➔  Choices are less noisy
➔  More affected by value
➔  More deterministic



SOLID  LIQUID  GAS

Tricia Seow | Samuel Hewitt | Noam Goldway

# How should we choose?
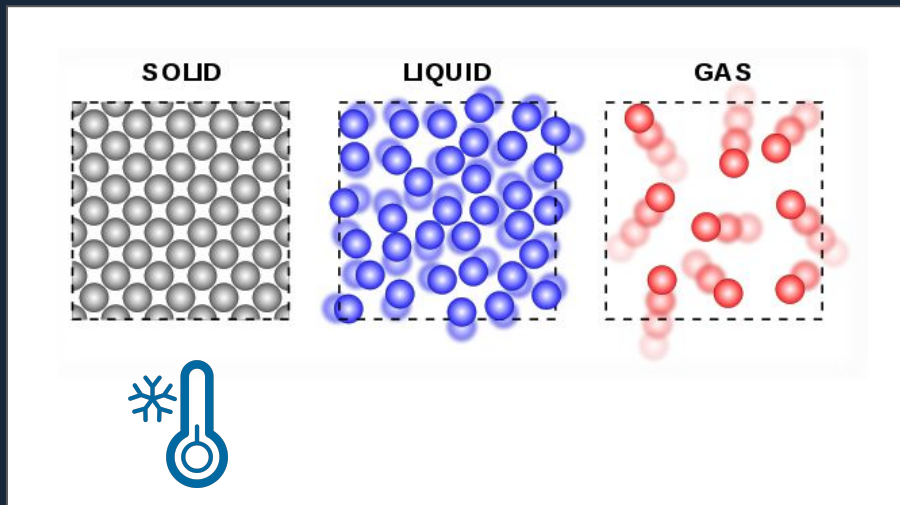
Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

## High temperature

➔ Choices are more noisy
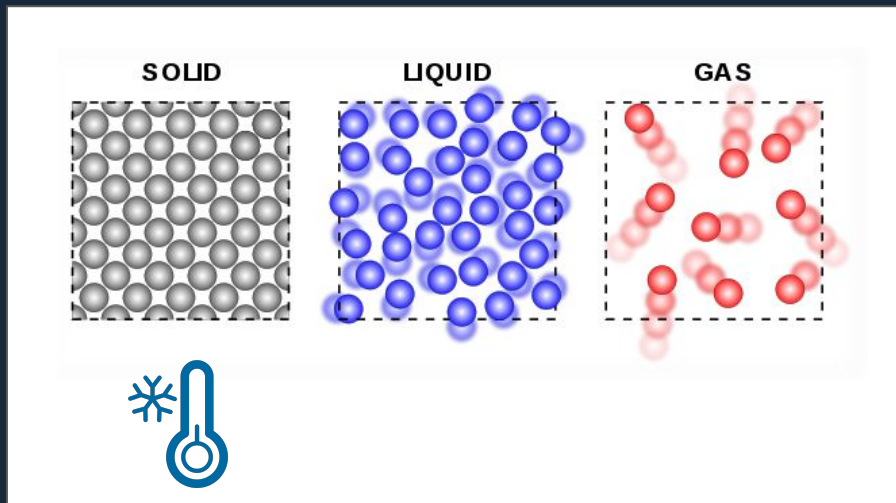➔ Less affected by value
➔ Less deterministic



SOLID     LIQUID     GAS

Tricia Seow | Samuel Hewitt | Noam Goldway

# How should we choose?

Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

→ Let's assume that if we don't pick **purple** we will pick **orange**; and vice versa

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

**Softmax equation:**

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\tau)}{\text{SUM}[\exp([\ ^{V_{\text{purple}}}_{V_{\text{orange}}}\ ]/\tau)]}$$

→ Let's assume that if we don't pick **purple** we will pick **orange**; and vice versa
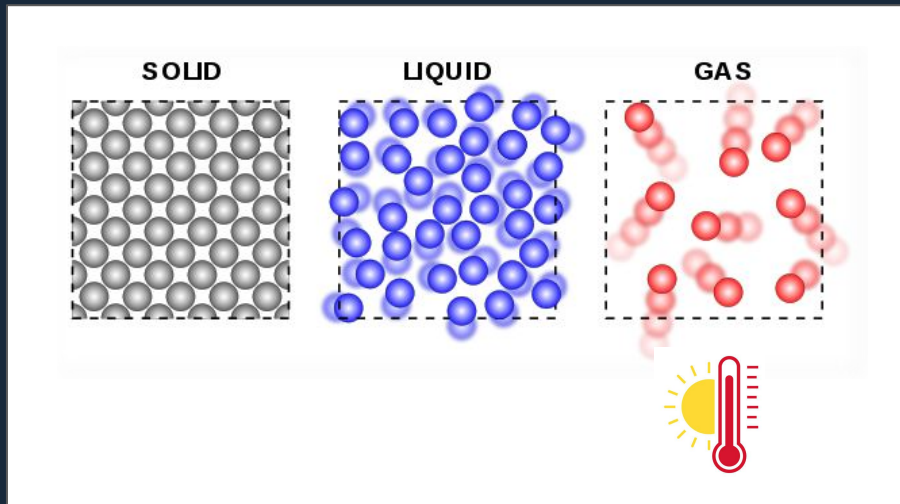
Tricia Seow | Samuel Hewitt | Noam Goldway

# How should we choose?

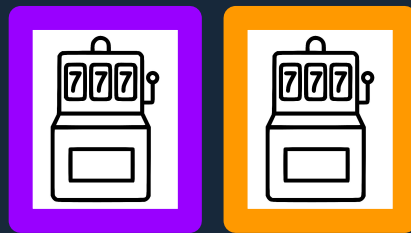Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

**Softmax equation:**

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\tau)}{SUM[\exp([\ V_{\text{purple}}\ V_{\text{orange}}\ ]/\tau)]}$$

Probability of choosing **purple**

→ $P(\text{orange}) = 1 - P(\text{purple})$

→ Let's assume that if we don't pick **purple** we will pick **orange**; and vice versa

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

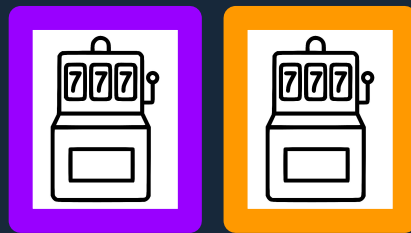**Softmax equation:**

Value of machines

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{SUM[\exp([\ \frac{V_{\text{purple}}}{V_{\text{orange}}}\ ]/\ \tau)]}$$

Probability of choosing **purple**

➔ $P(\text{orange}) = 1 - P(\text{purple})$

➔ Let's assume that if we don't pick **purple** we will pick **orange**; and vice versa

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?
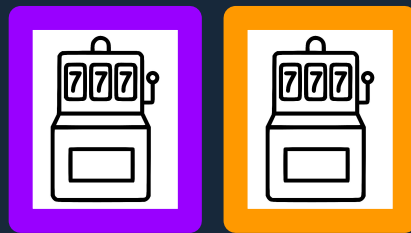
Temperature ($\tau$): parameter that determines the extent to which value estimates influence choice behaviour

**Softmax equation:**

Value of machines

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{\text{SUM}[\exp([\ \frac{V_{\text{purple}}}{V_{\text{orange}}}\ ]/\ \tau)]}$$

Probability of choosing **purple**

→  P(orange) = 1 - P(purple)

Free parameter temperature

→ Let's assume that if we don't pick **purple** we will pick **orange**; and vice versa

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{\text{orange}}^{\text{purple}}\ ]/\ \tau)]}$$

*Exploit*

➔ Choose slot machine when reward is better than the other

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

Choose **purple** is better ↑

Choose **orange** is better ↓



$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{\text{SUM}[\exp([V_{purple}\ \ V_{orange}\ ]/\ \tau)]}$$

*Exploit*

→ Choose slot machine when reward is better than the other

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(purple) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{SUM[\exp([\ V_{orange}^{purple}\ ]/\ \tau)]}$$

Choose **purple** is better

Choose **orange** is better



*Exploit*

➔ Choose slot machine when reward is better than the other

*Temperature is low*

➔ Choices are less noisy
➔ More affected by value
➔ More deterministic

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(\mathbf{purple}) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{SUM[\exp([\ V^{purple}_{orange}\ ]/\ \tau)]}$$

*Explore*

➔     Random choice

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

Choose **purple** is better

Choose **orange** is better

$$P(\textbf{purple}) = \frac{\exp([\ V_{\textbf{purple}}\ ]/\ \tau)}{SUM[\exp([\ V^{\textbf{purple}}_{\textbf{orange}}\ ]/\ \tau)]}$$

*Explore*

→   **Random choice**



Orange is better value

Purple is better value

tau = 10000

P(Choose Purple)

Value (Purple - Orange)

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB
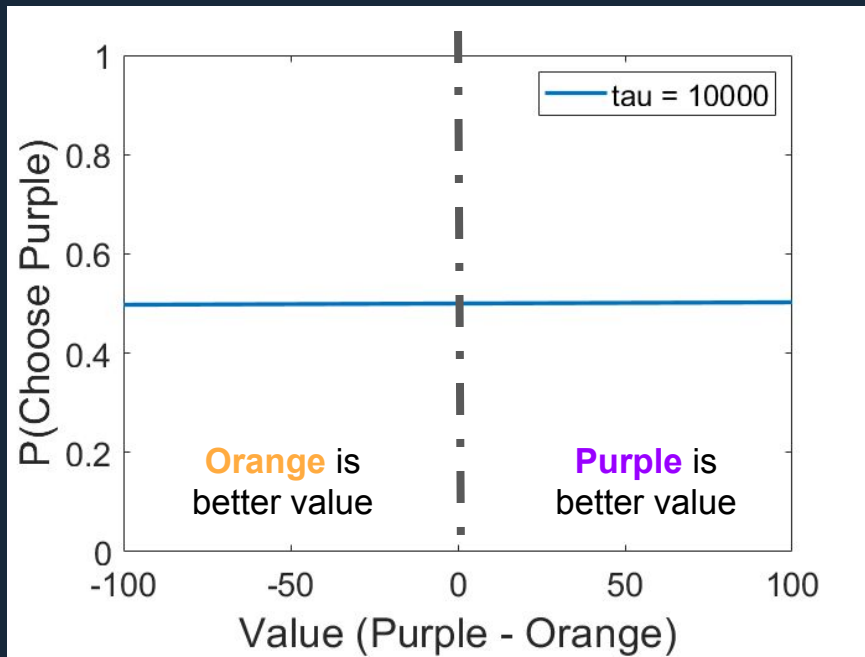
# How should we choose?

Choose **purple** is better

Choose **orange** is better



$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{SUM[\exp([\ V_{purple}^{purple}\ ]/\ \tau)]}$$

$V_{orange}$

*Explore*
➔ Random choice

*Temperature is high*
➔ Choices are more noisy
➔ Less affected by value
➔ More random

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{purple}\ ]/\ \tau)]}$$

$$V_{orange}$$

## *Temperature is low*

➔  Choices are less noisy
➔  More affected by value
➔  More deterministic

## *Temperature is high*

➔  Choices are more noisy
➔  Less affected by value
➔  Less deterministic

Tricia Seow | Samuel Hewitt | Noam Goldway

Developmental
Computational
Psychiatry lab

HARTLEY LAB

flux

# How should we choose?

Developmental
Computational
Psychiatry lab

HARTLEY LAB

Choose **purple** is better

Choose **orange** is better



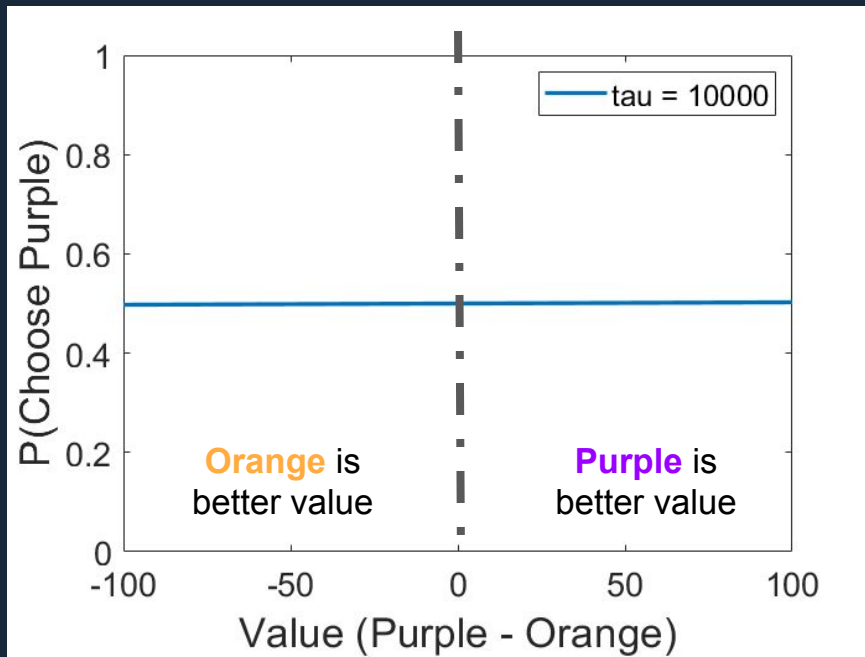$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{purple}\ ]/\ \tau)]}$$

$V_{orange}$

***Temperature is low***

➜ Choices are less noisy
➜ More affected by value
➜ More deterministic

***Temperature is high***

➜ Choices are more noisy
➜ Less affected by value
➜ Less deterministic

Tricia Seow | Samuel Hewitt | Noam Goldway

https://github.com/DevComPsy/2021FluxCompModellingWorkshop

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

*Softmax*

$$P(purple) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{SUM[\exp([\ V_{purple}^{purple}\ ]/\ \tau)]}$$

$V_{orange}$

**What does the exponential (exp) do?**

Tricia Seow | Samuel Hewitt | Noam Goldway

https://github.com/DevComPsy/2021FluxCompModellingWorkshop

flux

# How should we choose?

*Softmax*

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{orange}^{purple}\ ]/\ \tau)]}$$

**What does the exponential (exp) do?**

➔ **Non-linear transformation of value**

➔ **Deals with negative values**



Tricia Seow | Samuel Hewitt | Noam Goldway
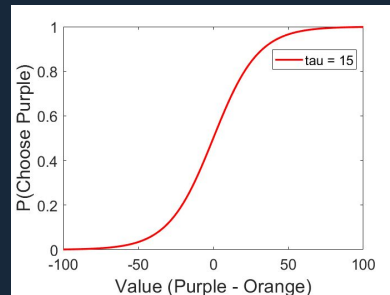
flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?



*Softmax*

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{\text{orange}}^{V_{\text{purple}}}\ ]/\ \tau)]}$$

**What does the exponential (exp) do?**

➔ Non-linear transformation of value
➔ Deals with negative values

**What does the division by SUM do?**

➔ Normalizes values to between 0 to 1

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ])}{\text{SUM}[\exp([\ V_{purple}^{purple}\ V_{orange}\ ])]}$$

*Softmax*

→ Transforms value input into values between 0 to 1

Assume temperature = 1

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ])}{SUM[\exp([\ V_{purple}\ V_{orange}\ ])]}$$

*Softmax*

→ Transforms value input into values between 0 to 1

Assume temperature = 1

For my next slot machine play...

$$V_{purple} = [\ 60\ ] \quad V_{orange} = [\ 40\ ]$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ])}{\text{SUM}[\exp([\ V_{purple},\ V_{orange}\ ])]}$$

$$= \frac{\exp([\ 60\ ])}{\text{SUM}[\exp([\ 60,\ 40\ ])]}$$

## *Softmax*

→ Transforms value input into values between 0 to 1

Assume temperature = 1

For my next slot machine play...

$$V_{purple} = [\ 60\ ] \quad V_{orange} = [\ 40\ ]$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ])}{\text{SUM}[\exp([V_{purple}\ V_{orange}\ ])]}$$

$$= \frac{\exp([\ 60\ ])}{\text{SUM}[\exp([\ 60\ 40\ ])]}$$

$$= \frac{e^{60}}{e^{60}+e^{40}}$$

$$= \quad 1$$

## *Softmax*

→ Transforms value input into values between 0 to 1
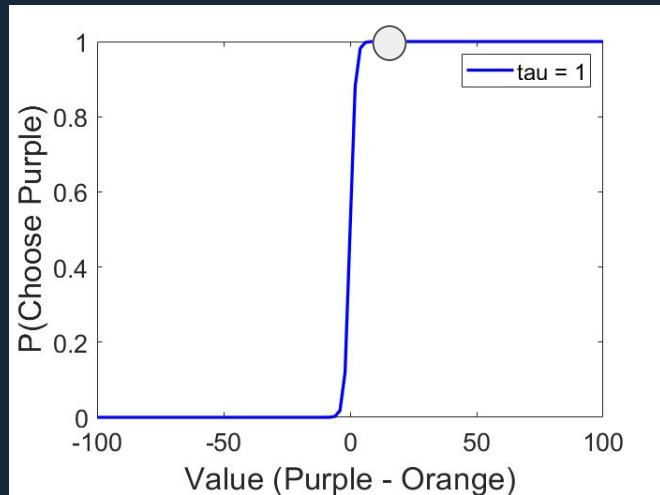
Assume temperature = 1

For my next slot machine play...

$$V_{purple} = [\ 60\ ] \quad V_{orange} = [\ 40\ ]$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

*Softmax*

➔ **Transforms value input into values between 0 to 1**

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ])}{\text{SUM}[\exp([\ V_{\text{purple}}\ V_{\text{orange}}\ ])]}$$

$$= \frac{\exp([\ 60\ ])}{\text{SUM}[\exp([\ 60\ 40\ ])]}$$

$$= \frac{e^{60}}{e^{60}+e^{40}}$$

$$= 1$$



$$V_{\text{purple}} = [\ 60\ ] \quad V_{\text{orange}} = [\ 40\ ]$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

*Softmax*

→ Transforms value input into values between 0 to 1

Assume temperature = 1

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ])}{SUM[\exp([\ V_{purple}\ V_{orange}\ ])]}$$

$$= \frac{\exp([\ 60\ ])}{SUM[\exp([\ 60\ 40\ ])]}$$

$$= \frac{e^{60}}{e^{60}+e^{40}}$$

$$= 1$$

$$P(\text{orange}) = \frac{\exp([\ V_{orange}\ ])}{SUM[\exp([\ V_{purple}\ V_{orange}\ ])]}$$

Tricia Seow | Samuel Hewitt | Noam Goldway

$$V_{purple} = [\ 60\ ] \quad V_{orange} = [\ 40\ ]$$

flux

# How should we choose?

*Softmax*

→ Transforms value input into values between 0 to 1

Assume temperature = 1

$$P(\text{purple}) = \frac{\exp([\ V_{purple}\ ])}{\text{SUM}[\exp([\ V_{purple}^{\ }\ V_{orange}\ ])]}$$

$$= \frac{\exp([\ 60\ ])}{\text{SUM}[\exp([\ 60\ 40\ ])]}$$

$$= \frac{e^{60}}{e^{60}+e^{40}}$$

$$= 1$$

$$P(\text{orange}) = \frac{\exp([\ V_{orange}\ ])}{\text{SUM}[\exp([\ V_{purple}^{\ }\ V_{orange}\ ])]}$$

$$= \frac{e^{40}}{e^{60}+e^{40}}$$

$$= 0$$

$$V_{purple} = [\ 60\ ] \quad V_{orange} = [\ 40\ ]$$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

## *Softmax*

→ Transforms value input into values between 0 to 1

Assume temperature = 1

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ])}{\text{SUM}[\exp([\ V_{\text{orange}}^{V_{\text{purple}}}\ ])]}$$

$$= \frac{\exp([\ 60\ ])}{\text{SUM}[\exp([\ \begin{smallmatrix}60\\40\end{smallmatrix}\ ])]}$$

$$P(\text{orange}) = \frac{\exp([\ V_{\text{orange}}\ ])}{\text{SUM}[\exp([\ V_{\text{orange}}^{V_{\text{purple}}}\ ])]}$$

$$= \frac{e^{60}}{e^{60}+e^{40}}$$

$$= \frac{e^{40}}{e^{60}+e^{40}}$$

$$= 1$$

$$= 0$$

probability equals to 1

Tricia Seow | Samuel Hewitt | Noam Goldway

$V_{\text{purple}} = [\ 60\ ]$  $V_{\text{orange}} = [\ 40\ ]$

flux

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{\text{orange}}^{\text{purple}}\ ]/\ \tau)]}$$

*Temperature ($\tau$)*

→ how much value affects choices

Assume temperature = 15

$V_{\text{purple}}$ = [ 60 ]   $V_{\text{orange}}$ = [ 40 ]

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{\text{SUM}[\exp([\ {}^{V_{\text{purple}}}_{V_{\text{orange}}}\ ]/\ \tau)]}$$

$$= \frac{\exp([\ 60\ ]/\ 15)}{\text{SUM}[\exp([\ {}^{60}_{40}\ ]/\ 15)]}$$

## *Temperature ($\tau$)*

→   how much value affects choices

Assume temperature = 15

$V_{\text{purple}}$ = [ 60 ]   $V_{\text{orange}}$ = [ 40 ]

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# How should we choose?

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\ \tau)}{\text{SUM}[\exp([\ V_{\text{orange}}^{\text{purple}}\ ]/\ \tau)]}$$

$$= \frac{\exp([\ 60\ ]/\ 15)}{\text{SUM}[\exp([\ \begin{matrix} 60 \\ 40 \end{matrix}\ ]/\ 15)]}$$

$$= \frac{e^{60/15}}{e^{60/15}+e^{40/15}}$$

$$= \quad 0.79$$

## Temperature ($\tau$)

→  how much value affects choices

Assume temperature = 15

$V_{\text{purple}} = [\ 60\ ]$  $V_{\text{orange}} = [\ 40\ ]$

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB

# How should we choose?

## Temperature ($\tau$)

→ how much value affects choices

Assume temperature = 15

$$P(\text{purple}) = \frac{\exp([\text{ V}_{\text{purple}}\text{ ]}/\tau)}{\text{SUM}[\exp([\text{ V}_{\text{orange}}^{\text{purple}}\text{ ]}/\tau)]}$$

$$= \frac{\exp([\text{ 60 ]}/15)}{\text{SUM}[\exp([\overset{60}{\underset{40}{}}\text{ ]}/15)]}$$

$$P(\text{orange}) = \frac{\exp([\text{ V}_{\text{orange}}\text{ ]}/\tau)}{\text{SUM}[\exp([\text{ V}_{\text{orange}}^{\text{purple}}\text{ ]}/\tau)]}$$

$$= \frac{e^{60/15}}{e^{60/15}+e^{40/15}}$$

$$= \frac{e^{40/15}}{e^{60/15}+e^{40/15}}$$

$$= \quad 0.79 \qquad\qquad\qquad = \quad 0.21$$

probability equals to 1

Tricia Seow | Samuel Hewitt | Noam Goldway

$$V_{\text{purple}} = [\,60\,] \quad V_{\text{orange}} = [\,40\,]$$
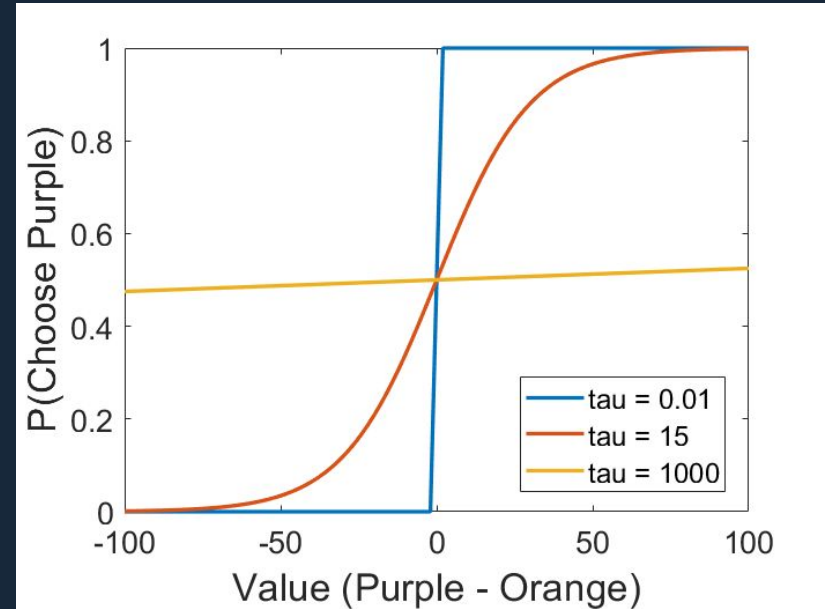
flux

# What have we learnt about choice?

**Softmax equation:**

$$P(\text{purple}) = \frac{\exp([\ V_{\text{purple}}\ ]/\tau)}{SUM[\exp([\ V_{\text{purple}}\ ]/\tau)]}$$

$$V_{\text{orange}}$$

**Temperature ($\tau$):** parameter that determines the extent to which value estimates influence choice behaviour

**Exploit or Explore**



Tricia Seow | Samuel Hewitt | Noam Goldway

flux

Developmental
Computational
Psychiatry lab

HARTLEY LAB
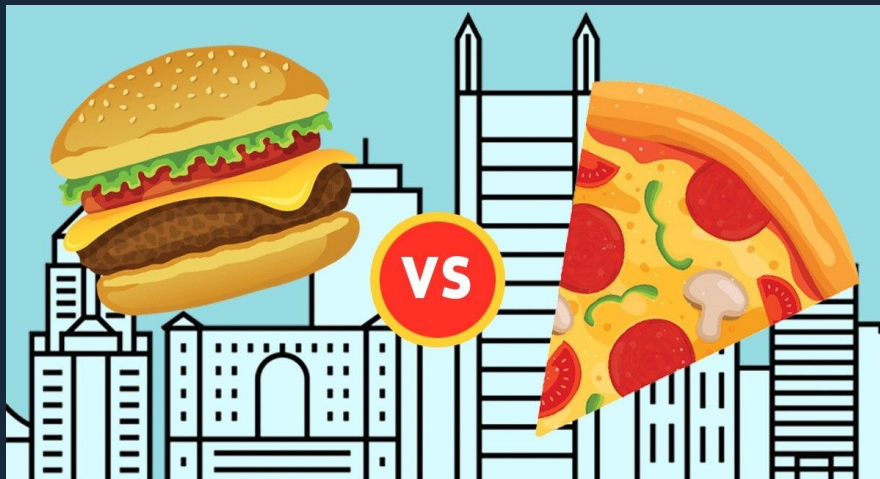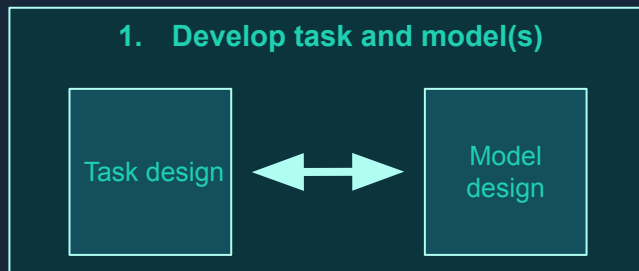
# What have we learnt about choice?

## Exploit versus Explore:

➔ Discover "what works" by alternating between exploration and exploitation

➔ In uncertain environments, more exploration could be useful



VS

Tricia Seow | Samuel Hewitt | Noam Goldway

flux

# Summary

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

Task design ⟷ Model design

| Task | Trial and error learning | Reinforcement Learning |
|---|---|---|
| Value Function | Subjective value from objective outcomes | Rescorla-Wagner |
| Choice Function | Choice probabilities from value | Softmax |

Tricia Seow | Samuel Hewitt | Noam Goldway

flux