

Section B | 30 minutes

How to Develop a Computational Model?

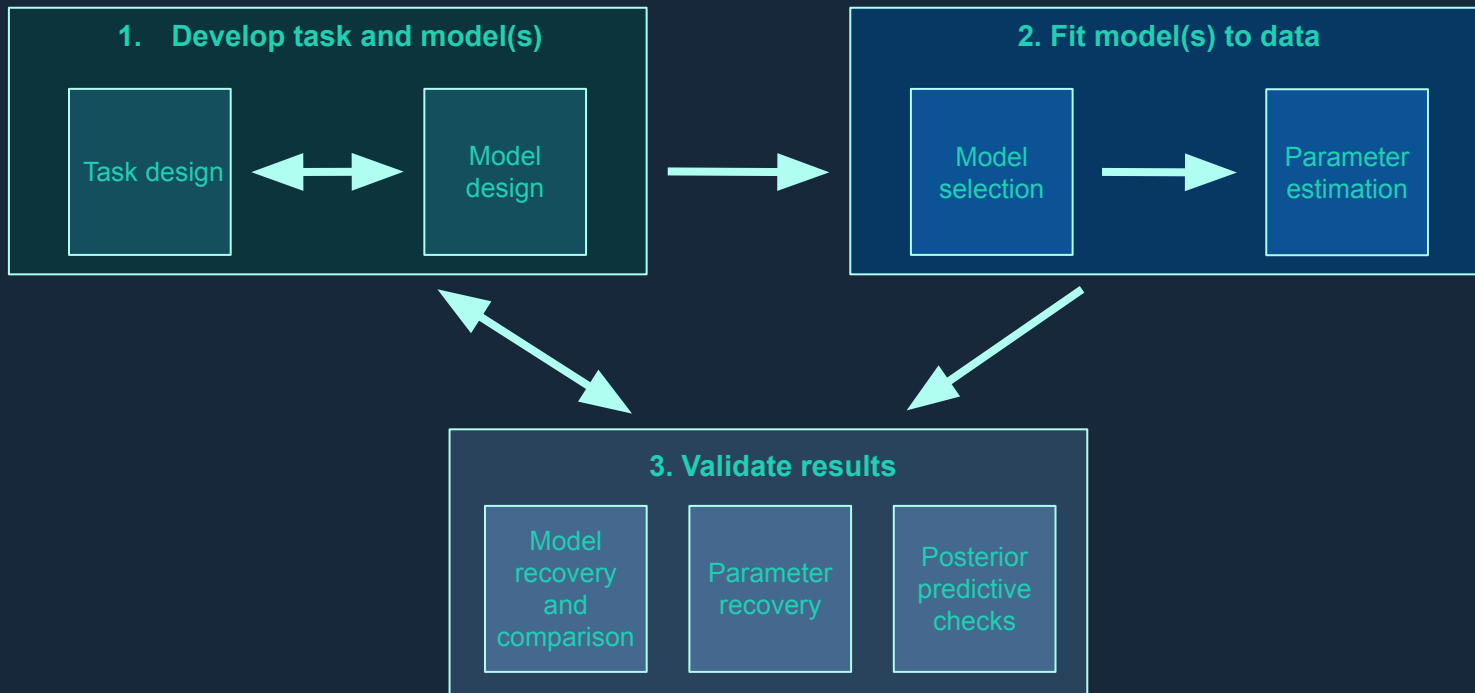
"All models are wrong, but some are useful"
George E. P. Box

Tricia Seow | Samuel Hewitt | Noam Goldway

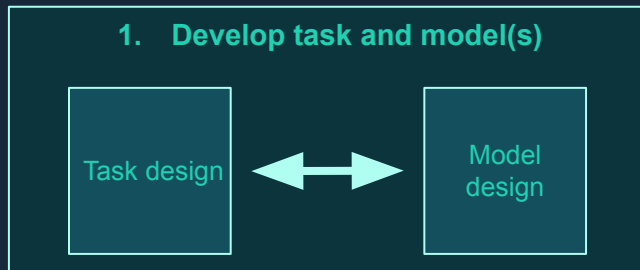
What we will cover:

- 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

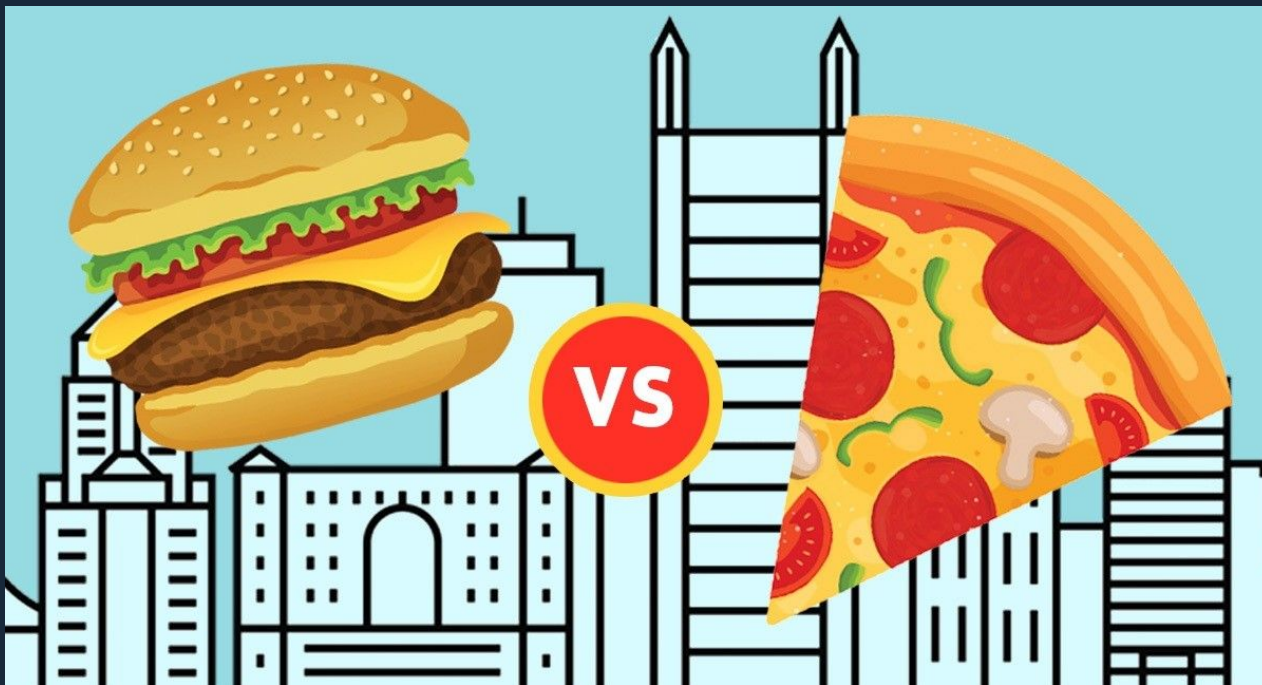
Developing a computational model



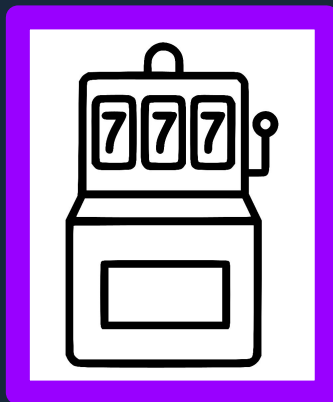
Developing a computational model



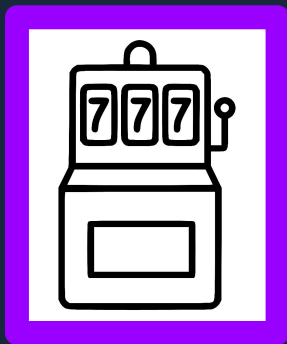
Section B | How to Develop a Computational Model?



2-arm bandit



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

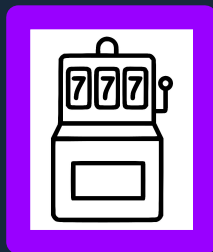
Modelling behaviour with RL



Value function: Rescorla Wagner model

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

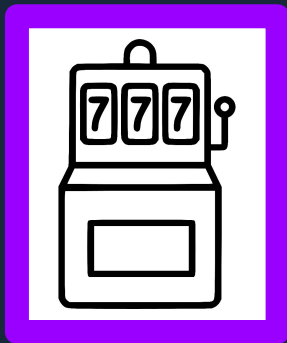
Prediction Error



$V_{\text{prepule}} > V_{\text{orange}}$



Prediction Error



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

Prediction Error



?

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

Prediction Error



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

Prediction Error



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

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)

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

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

(Reward - Value on
previous trial)

Prediction error
what you received - what you expected

Modelling behaviour with RL

Value function: Rescorla Wagner model

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

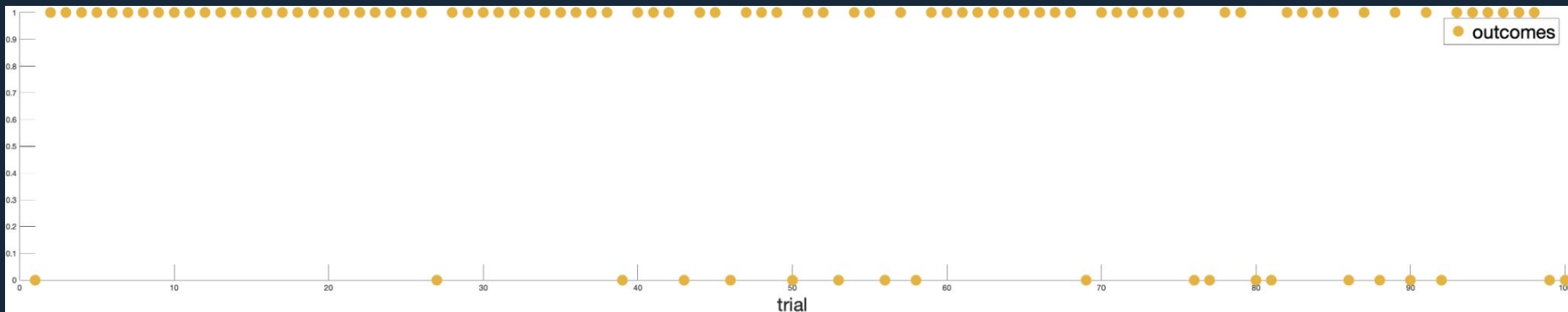
Value
(of the slot machine) = Value on previous trial + Learning rate (Reward - Value on previous trial)

Prediction error
what you received - what you expected

Modelling behaviour with RL

Prediction error
what you received - what you expected

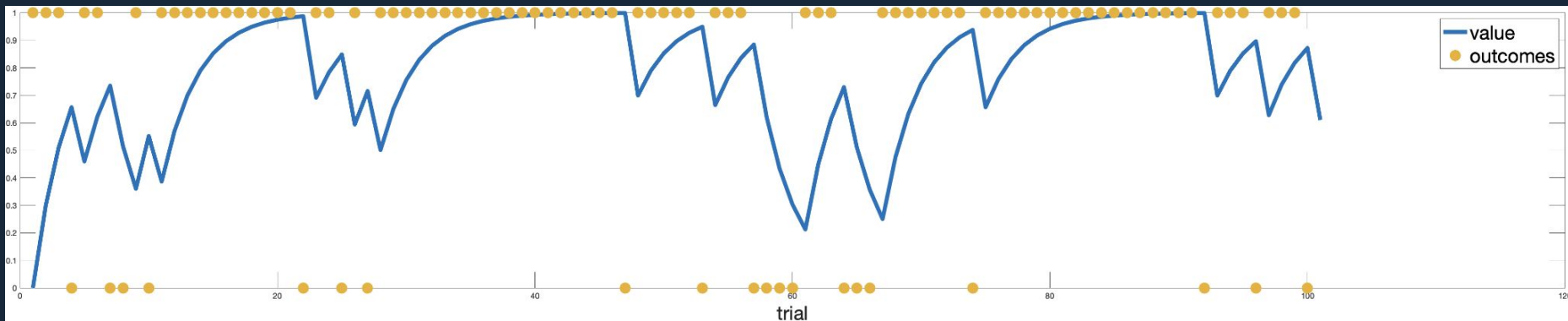
$$\text{Value (of the slot machine)} = \text{Value on previous trial} + \text{Learning rate} \left(\text{Reward} - \text{Value on previous trial} \right)$$



Modelling behaviour with RL

Prediction error
what you received - what you expected

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



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \underbrace{\left(\begin{array}{c} \text{Prediction error} \\ \text{what you received - what you expected} \\ \text{Reward - Value on} \\ \text{previous trial} \end{array} \right)}$$

<i>Trial</i> 1	?				
-------------------	---	--	--	--	--



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \underbrace{\left(\begin{array}{c} \text{Reward} \\ \text{Prediction error} \\ \text{what you received - what you expected} \end{array} - \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} \right)}$$

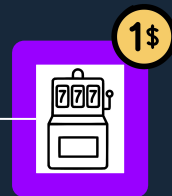
<i>Trial</i> 1	?	<i>initiation: 0.5</i>			<i>initiation: 0.5</i>
-------------------	---	------------------------	--	--	------------------------



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \left(\begin{array}{c} \text{Reward} \\ \text{Prediction error} \\ \text{what you received - what you expected} \end{array} - \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} \right)$$

<i>Trial</i> 1	?	0.5		1	0.5
-------------------	---	-----	--	---	-----



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \underbrace{\left(\begin{array}{c} \text{Reward} \\ \text{Prediction error} \\ \text{what you received - what you expected} \end{array} - \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} \right)}$$

<i>Trial</i> 1	?	0.5	+	1	0.5
-------------------	---	-----	---	---	-----



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \left(\begin{array}{c} \text{Prediction error} \\ \text{what you received - what you expected} \\ \text{Reward} - \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} \end{array} \right)$$

<i>Trial</i> 1	1	0.5	+	0.5
-------------------	---	-----	---	-----

Modelling behaviour with RL

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

Prediction error
what you received - what you expected

Trial 1	1	0.5	+	0.5
Trial 2		1	+	1

Modelling behaviour with RL

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

Prediction error
what you received - what you expected

Trial 1	1	0.5	+	0.5
Trial 2		1	+	1



Modelling behaviour with RL

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

Prediction error
what you received - what you expected

Trial 1	1	0.5	+	0.5
Trial 2		1	+	(1 - 1)



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \left(\begin{array}{c} \text{Prediction error} \\ \text{what you received - what you expected} \\ \text{Reward - Value on} \\ \text{previous trial} \end{array} \right)$$

<i>Trial</i> 1	1	0.5	+	0.5
<i>Trial</i> 2		1	+	0

Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \left(\begin{array}{c} \text{Prediction error} \\ \text{what you received - what you expected} \\ \text{Reward - Value on} \\ \text{previous trial} \end{array} \right)$$

<i>Trial</i> 1	1	0.5	+	0.5
<i>Trial</i> 2	1	1	+	0

Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \left(\begin{array}{c} \text{Prediction error} \\ \text{what you received - what you expected} \\ \text{Reward - Value on} \\ \text{previous trial} \end{array} \right)$$

<i>Trial</i> 1	1	0.5	+	0.5
<i>Trial</i> 2	1	1	+	0

Modelling behaviour with RL

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

Prediction error
what you received - what you expected

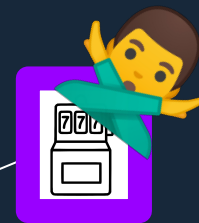
Trial 1	1	0.5	+	0.5
Trial 2	1	1	+	0
Trial 3		1	+	1

Modelling behaviour with RL

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

Prediction error
what you received - what you expected

Trial 1	1	0.5	+	0.5
Trial 2	1	1	+	0
Trial 3		1	+	0 1



Modelling behaviour with RL

$$\begin{array}{c} \text{Value} \\ \text{(of the slot machine)} \end{array} = \begin{array}{c} \text{Value on} \\ \text{previous trial} \end{array} + \begin{array}{c} \text{Learning} \\ \text{rate} \end{array} \underbrace{\left(\begin{array}{c} \text{Prediction error} \\ \text{what you received - what you expected} \\ \text{Reward - Value on} \\ \text{previous trial} \end{array} \right)}$$

<i>Trial</i> 1	1	0.5	+	0.5
<i>Trial</i> 2	1	1	+	0
<i>Trial</i> 3		1	+	-1



Modelling behaviour with RL

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

Prediction error
what you received - what you expected

<i>Trial</i> 1	1	0.5	+	0.5
<i>Trial</i> 2	1	1	+	0
<i>Trial</i> 3	0	1	+	-1

Modelling behaviour with RL

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

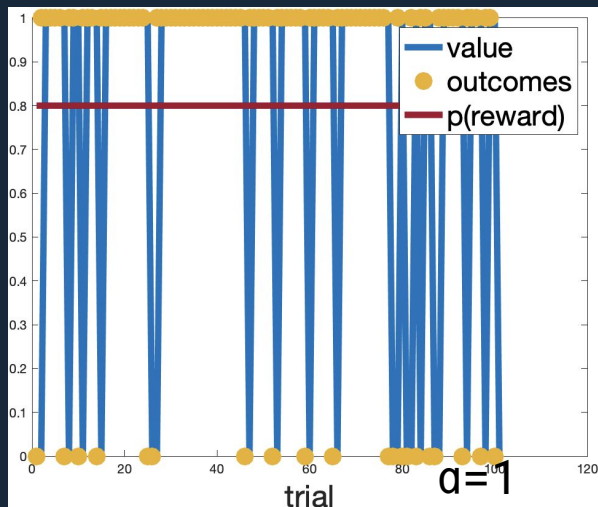
Prediction error
what you received - what you expected

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

How much should we learn?

What happens if we manipulate learning rate?

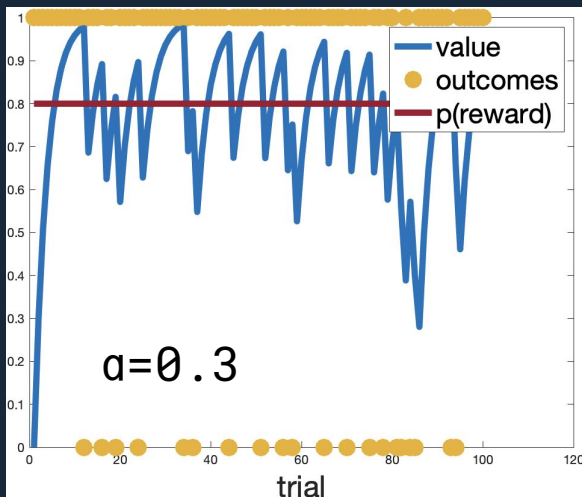
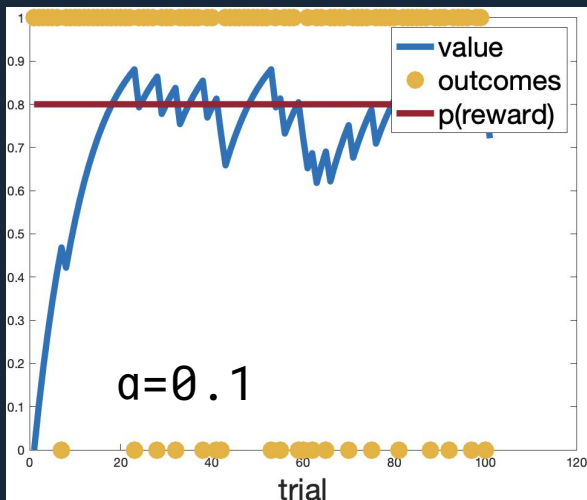
$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$



How much should we learn?

What happens if we manipulate learning rate?

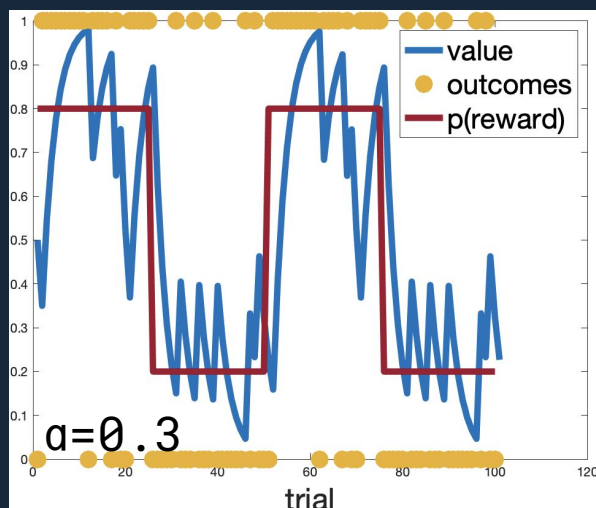
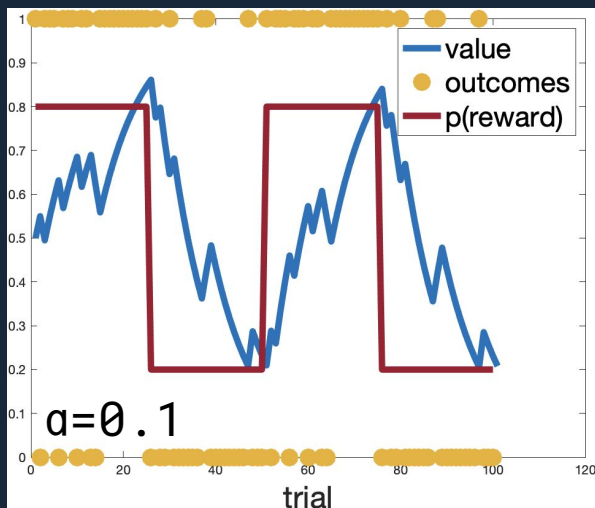
$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$



How much should we learn?

What happens if we manipulate learning rate?

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



Is low learning rate always better?

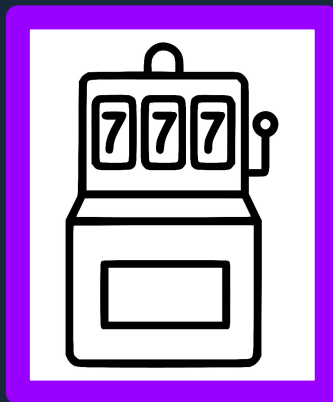
$$V_t = V_{t-1} + \alpha(R_t - V_{t-1})$$

- Depend on the statistics of the environment
- Low volatility- \rightarrow low α is better
 - High volatility- \rightarrow high α is better

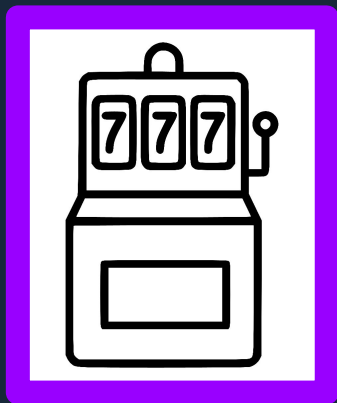
What did we learn so far

- 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

How should we choose?



How should we choose?



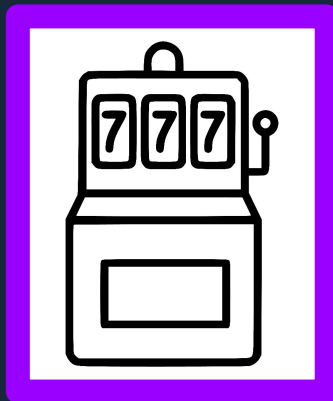
80% reward



20% reward

← learnt via trial and error
(value function) →

How should we choose?

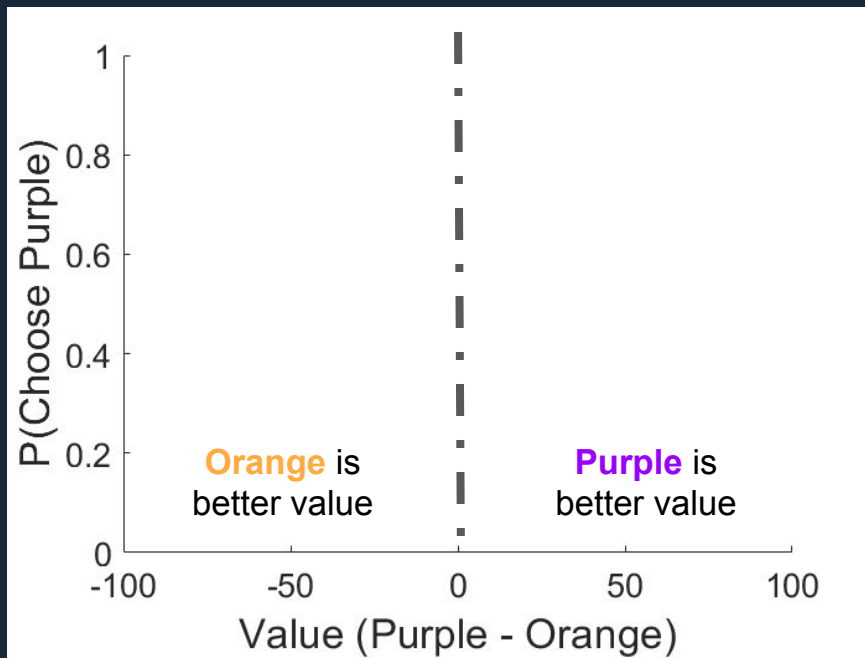


80% reward

Maximise rewards

- Pick slot machine with largest likelihood of reward
- Exploit

How should we choose?

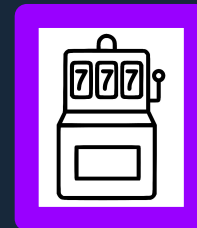
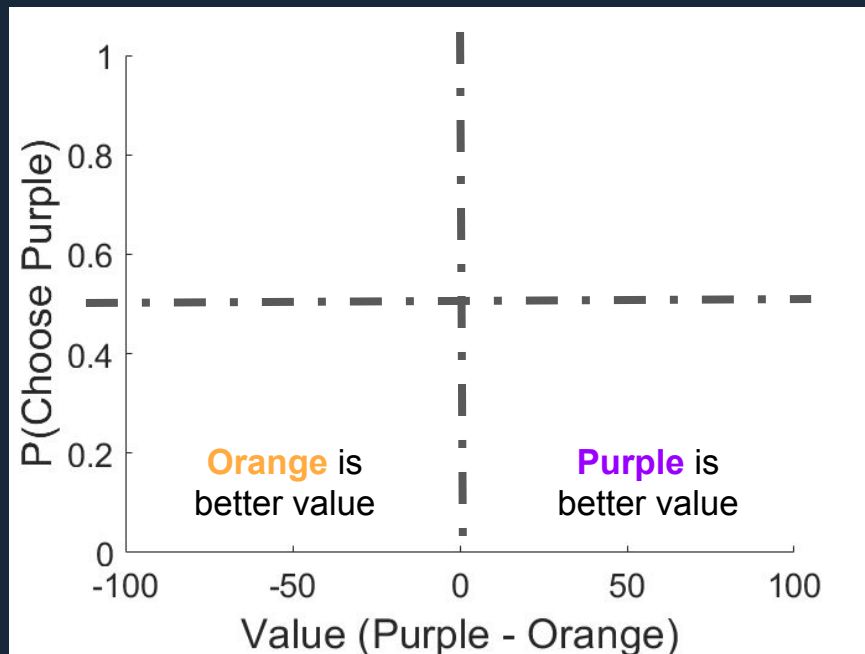


How should we choose?

Choose
purple is
better



Choose
orange is
better

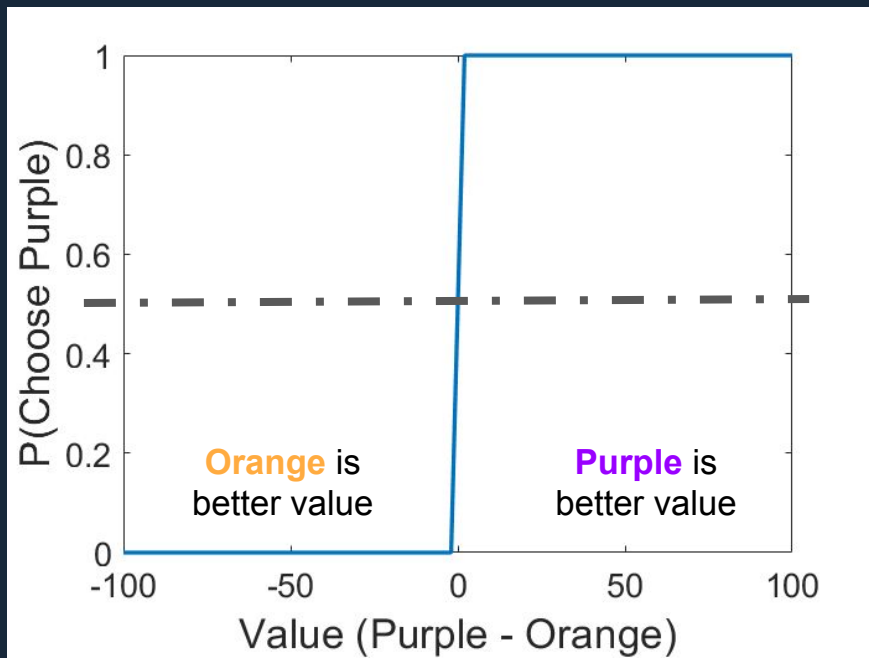


How should we choose?

Choose
purple is
better



Choose
orange is
better



Exploit

→ Choose slot machine when
reward is better than the
other

How should we choose?

Try other options

- Sample the outcomes of the other slot machine
- Explore



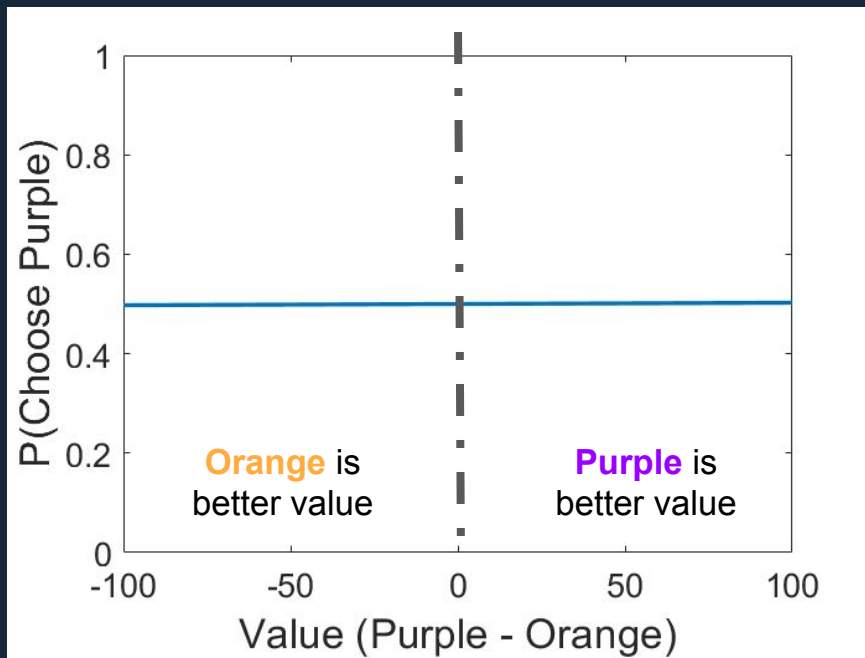
20% reward

How should we choose?

Choose
purple is
better



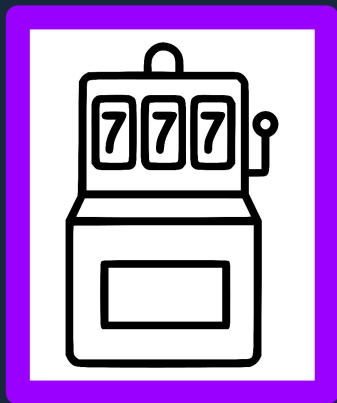
Choose
orange is
better



Explore

→ Choose slot machine
equally

How should we choose?



Exploit → an individual difference we can model as a free parameter ← *Explore*

How should we choose?

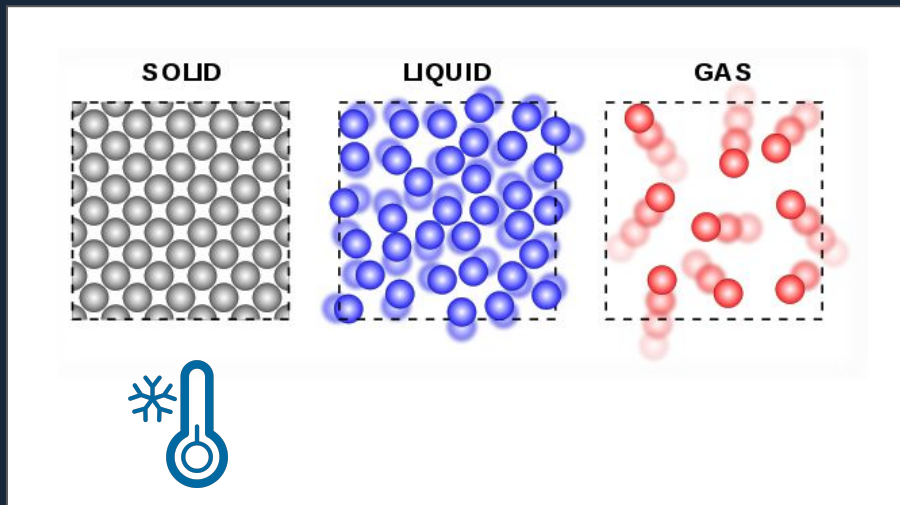
Temperature (τ): parameter that determines the extent to which value estimates influence choice behaviour

How should we choose?

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

Low temperature

→ Choices are less noisy

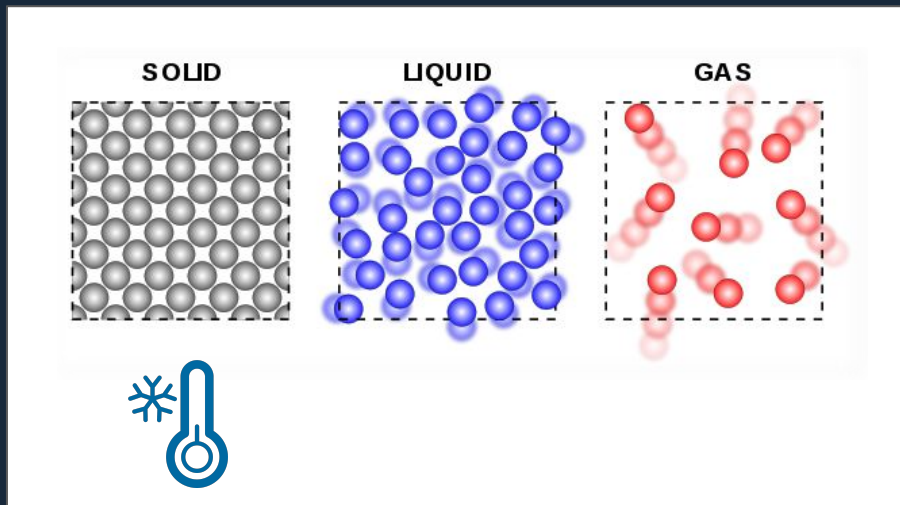


How should we choose?

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

Low temperature

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

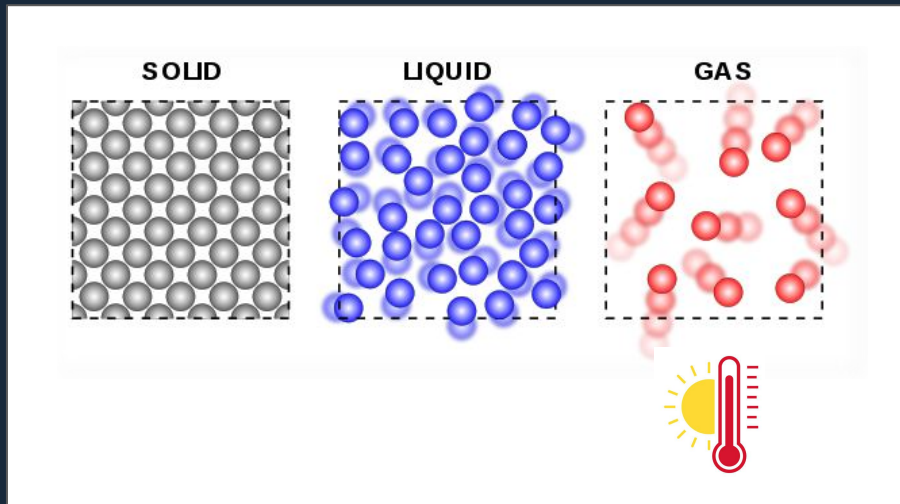


How should we choose?

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

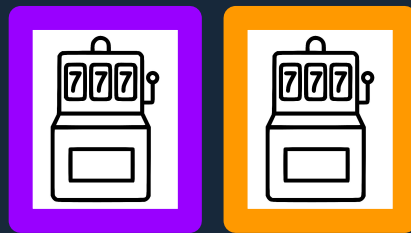
High temperature

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



How should we choose?

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



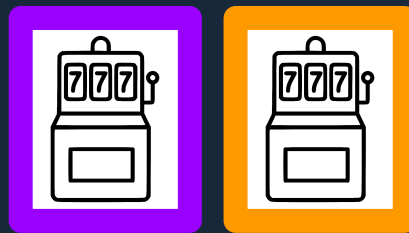
→ Let's assume that if we don't pick **purple** we will pick **orange**; and vice versa

How should we choose?

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

Softmax equation:

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



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

How should we choose?

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

Softmax equation:

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

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



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

How should we choose?

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

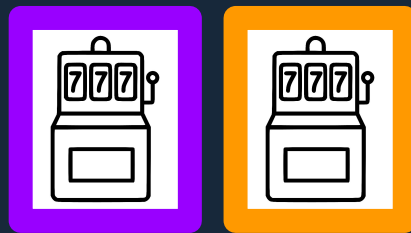
Softmax equation:

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

Value of machines

Probability of choosing purple

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



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

How should we choose?

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

Softmax equation:

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

Value of machines

Probability of choosing purple

Free parameter temperature

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



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

How should we choose?

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

Exploit

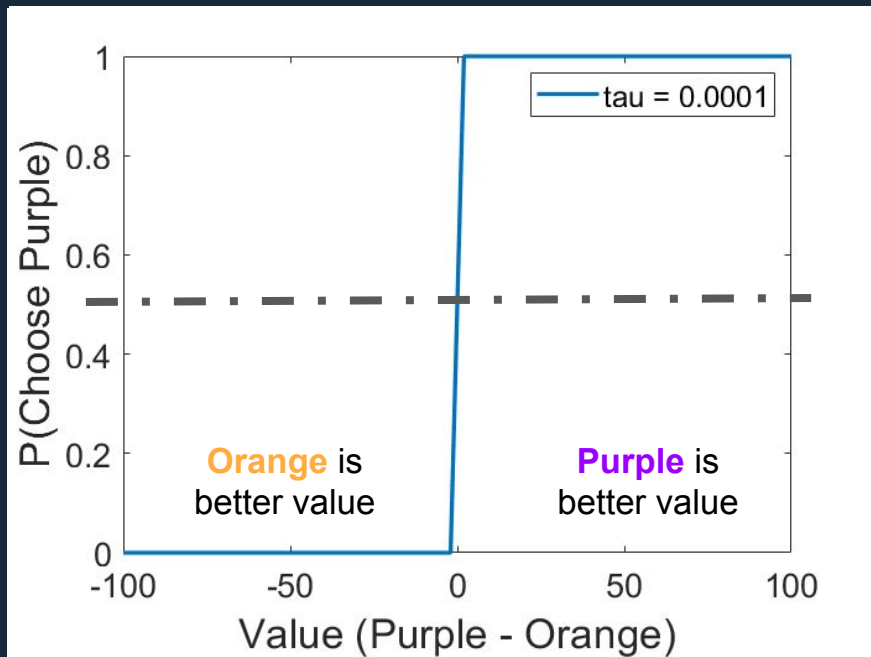
- Choose slot machine when reward is better than the other

How should we choose?

Choose
purple is
better



Choose
orange is
better



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

Exploit

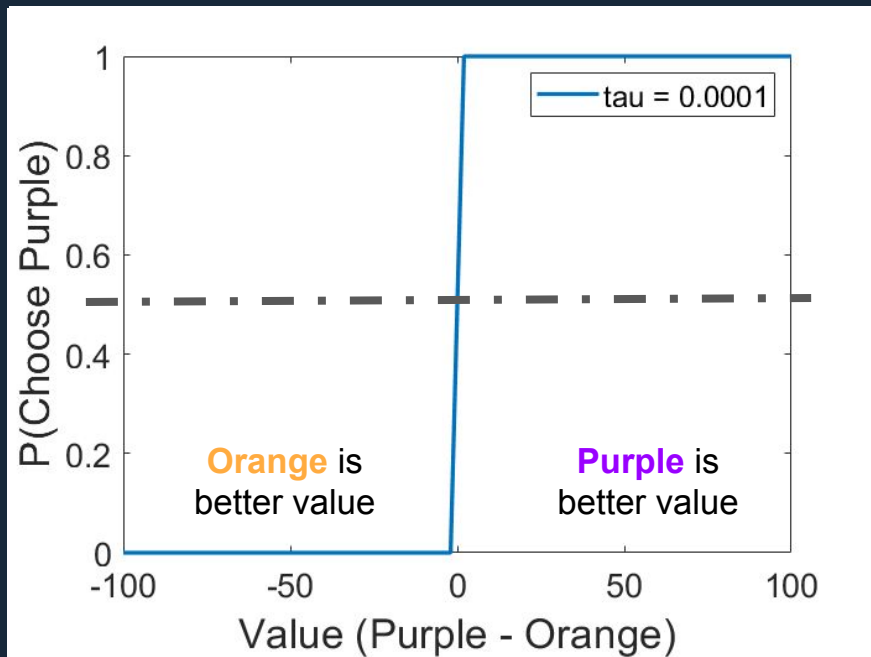
→ Choose slot machine when
reward is better than the
other

How should we choose?

Choose
purple is
better



Choose
orange is
better



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

Exploit

→ Choose slot machine when reward is better than the other

Temperature is low

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



How should we choose?

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

Explore

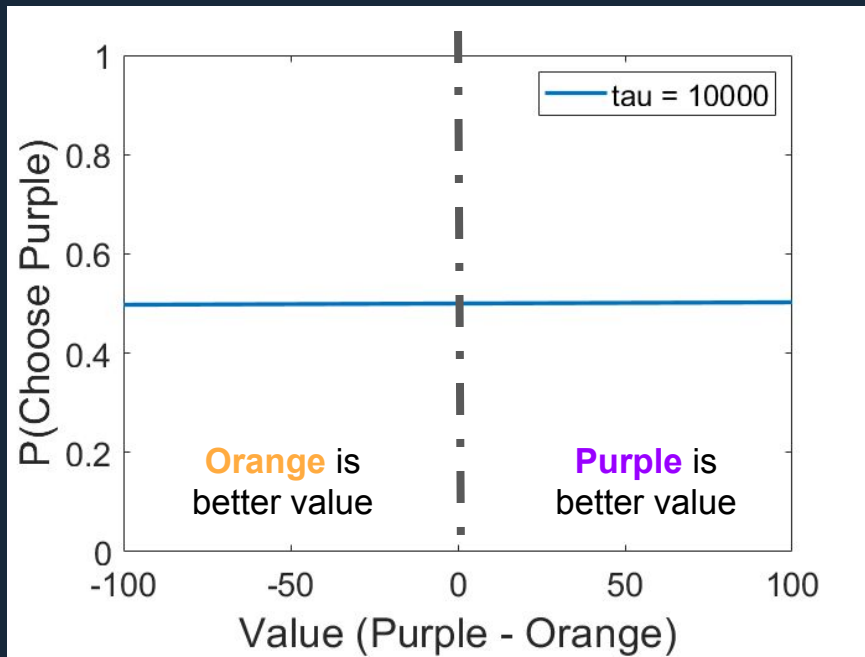
→ Random choice

How should we choose?

Choose
purple is
better



Choose
orange is
better



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

Explore

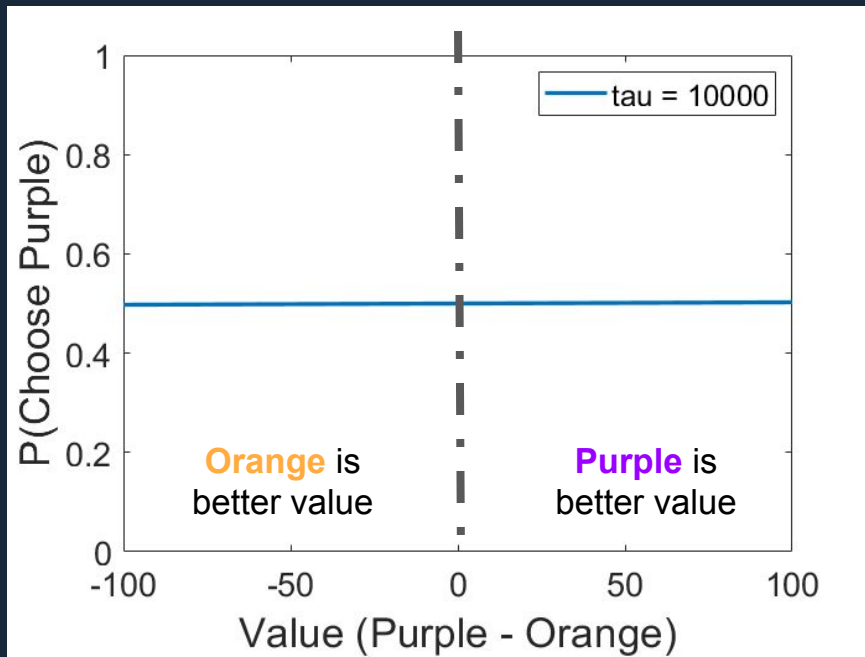
→ Random choice

How should we choose?

Choose
purple is
better



Choose
orange is
better



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

Explore

→ Random choice

Temperature is high

- Choices are more noisy
- Less affected by value
- More random



How should we choose?

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

Temperature is low

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

Temperature is high

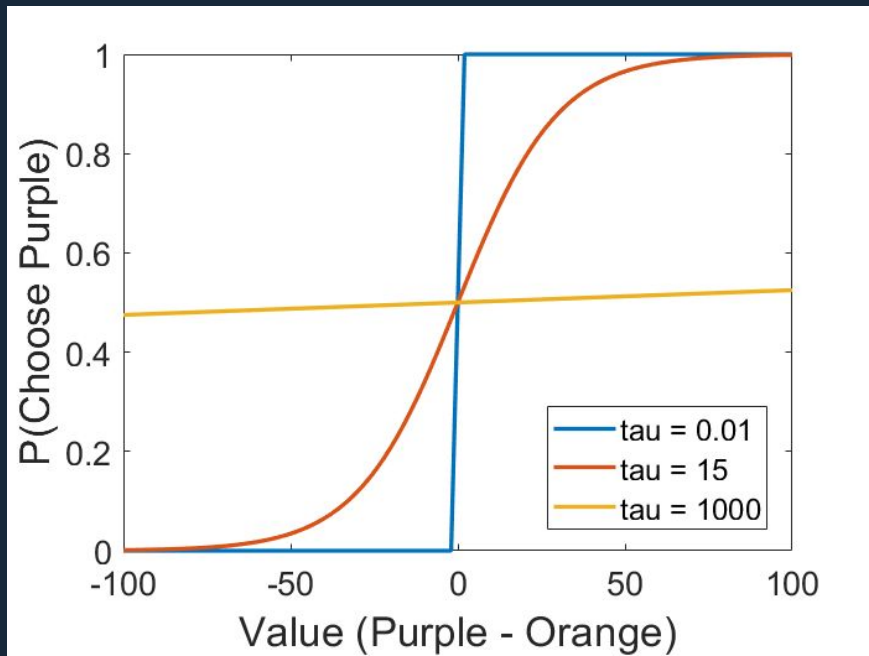
- Choices are more noisy
- Less affected by value
- Less deterministic

How should we choose?

Choose
purple is
better



Choose
orange is
better



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

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

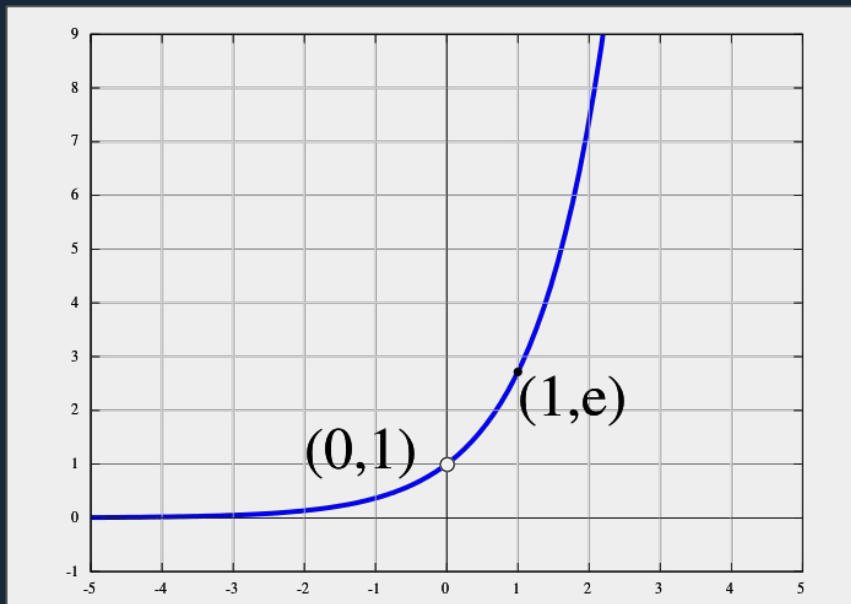
How should we choose?

Softmax

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

What does the exponential (exp) do?

How should we choose?

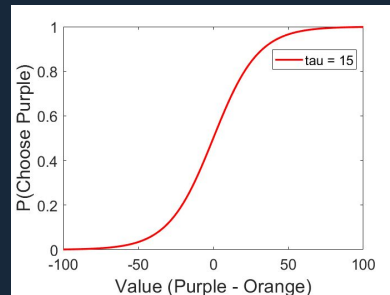


Softmax

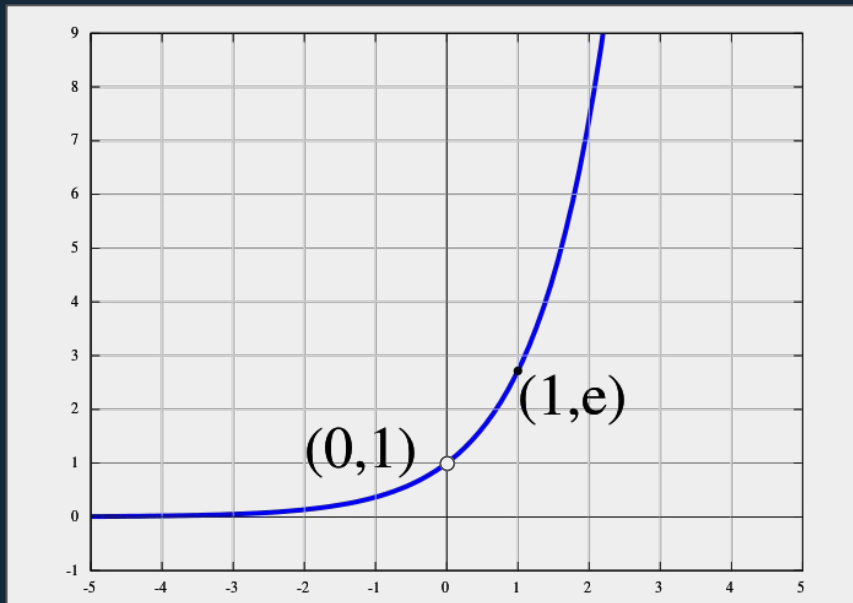
$$P(\text{purple}) = \frac{\exp([V_{\text{purple}}] / \tau)}{\text{SUM}[\exp([V_{\text{purple}}] / \tau)]}$$

What does the exponential (exp) do?

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



How should we choose?



Softmax

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

What does the exponential (exp) do?

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

What does the division by SUM do?

- Normalizes values to between 0 to 1

How should we choose?

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

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

How should we choose?

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

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

For my next slot machine play...

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

How should we choose?

$$\begin{aligned} P(\text{purple}) &= \frac{\exp([V_{\text{purple}}])}{\text{SUM}[\exp([V_{\text{purple}}], V_{\text{orange}}])} \\ &= \frac{\exp([60])}{\text{SUM}[\exp([60], 40)]} \end{aligned}$$

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

For my next slot machine play...

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

How should we choose?

$$\begin{aligned} P(\text{purple}) &= \frac{\exp([V_{\text{purple}}])}{\text{SUM}[\exp([V_{\text{purple}}], V_{\text{orange}}])} \\ &= \frac{\exp([60])}{\text{SUM}[\exp([60], 40)]} \\ &= \frac{e^{60}}{e^{60} + e^{40}} \\ &= 1 \end{aligned}$$

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

For my next slot machine play...

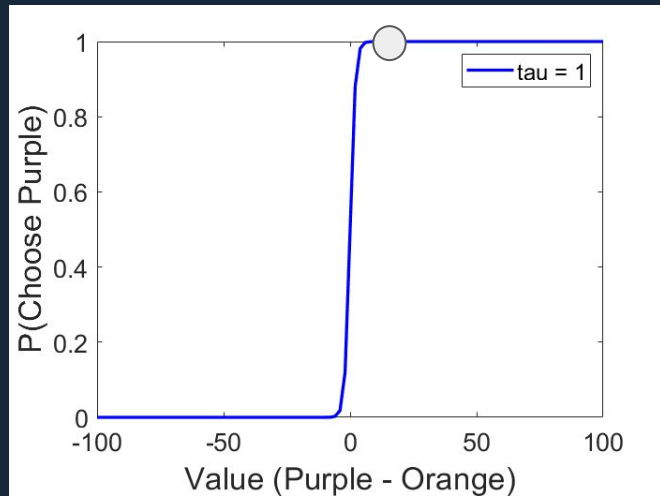
$$V_{\text{purple}} = [60] \quad V_{\text{orange}} = [40]$$

How should we choose?

$$\begin{aligned}
 P(\text{purple}) &= \frac{\exp([V_{\text{purple}}])}{\text{SUM}[\exp([V_{\text{purple}}], [V_{\text{orange}}])]} \\
 &= \frac{\exp([60])}{\text{SUM}[\exp([60], [40])]} \\
 &= \frac{e^{60}}{e^{60} + e^{40}} \\
 &= 1
 \end{aligned}$$

Softmax

→ Transforms value input into values between 0 to 1



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

How should we choose?

$$\begin{aligned} P(\text{purple}) &= \frac{\exp([V_{\text{purple}}])}{\text{SUM}[\exp([V_{\text{purple}}], [V_{\text{orange}}])]} \\ &= \frac{\exp([60])}{\text{SUM}[\exp([60], [40])]} \\ &= \frac{e^{60}}{e^{60} + e^{40}} \\ &= 1 \end{aligned}$$

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

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

How should we choose?

$$\begin{aligned}
 P(\text{purple}) &= \frac{\exp([V_{\text{purple}}])}{\text{SUM}[\exp([V_{\text{purple}}], V_{\text{orange}}])} \\
 &= \frac{\exp([60])}{\text{SUM}[\exp([60], 40])} \\
 &= \frac{e^{60}}{e^{60} + e^{40}} \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 P(\text{orange}) &= \frac{\exp([V_{\text{orange}}])}{\text{SUM}[\exp([V_{\text{purple}}], V_{\text{orange}}])} \\
 &= \frac{e^{40}}{e^{60} + e^{40}} \\
 &= 0
 \end{aligned}$$

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

How should we choose?

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

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

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

$$= 1$$

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

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

$$= 0$$

probability equals to 1

Softmax

→ Transforms value input into values between 0 to 1

Assume temperature = 1

How should we choose?

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

Temperature (τ)

→ how much value affects choices

Assume temperature = 15

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

How should we choose?

$$\begin{aligned} P(\text{purple}) &= \frac{\exp([V_{\text{purple}}] / \tau)}{\text{SUM}[\exp([V_{\text{purple}}] / \tau)]} \\ &= \frac{\exp([60] / 15)}{\text{SUM}[\exp([60] / 15)]} \end{aligned}$$

Temperature (τ)

→ how much value affects choices

Assume temperature = 15

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

How should we choose?

$$\begin{aligned} P(\text{purple}) &= \frac{\exp([V_{\text{purple}}] / \tau)}{\text{SUM}[\exp([V_{\text{purple}}] / \tau)]} \\ &= \frac{\exp([60] / 15)}{\text{SUM}[\exp([60] / 15)]} \\ &= \frac{e^{60/15}}{e^{60/15} + e^{40/15}} \\ &= 0.79 \end{aligned}$$

Temperature (τ)

→ how much value affects choices

Assume temperature = 15

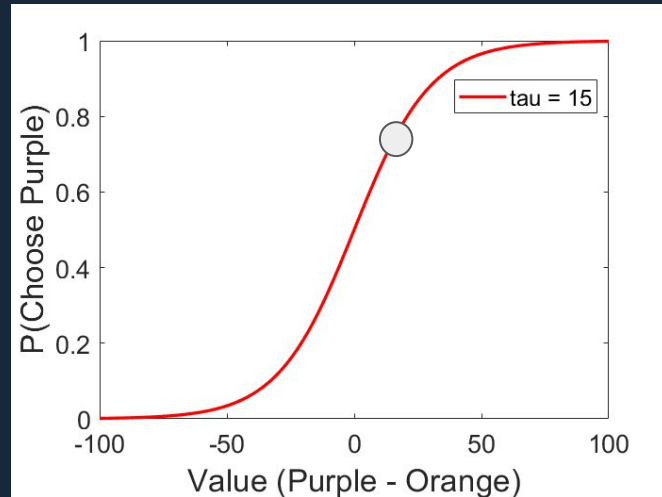
$$V_{\text{purple}} = [60] \quad V_{\text{orange}} = [40]$$

How should we choose?

$$\begin{aligned}
 P(\text{purple}) &= \frac{\exp([V_{\text{purple}}] / \tau)}{\text{SUM}[\exp([V_{\text{purple}}] / \tau)]} \\
 &= \frac{\exp([60] / 15)}{\text{SUM}[\exp([60] / 15)]} \\
 &= \frac{e^{60/15}}{e^{60/15} + e^{40/15}} \\
 &= 0.79
 \end{aligned}$$

Temperature (τ)

→ how much value affects choices



How should we choose?

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

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

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

$$= 0.79$$

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

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

$$= 0.21$$

probability equals to 1

Temperature (τ)

→ how much value affects choices

Assume temperature = 15

What have we learnt about choice?

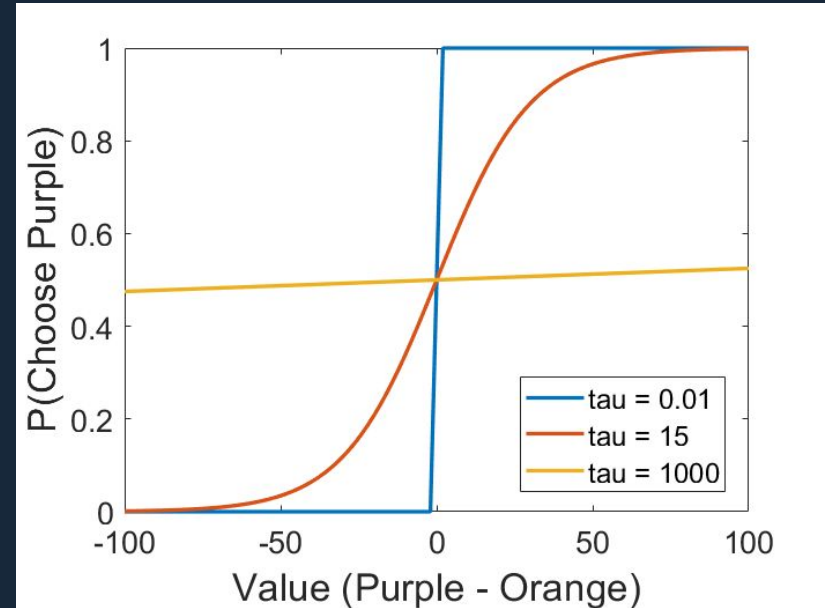
Softmax equation:

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

V_{orange}

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

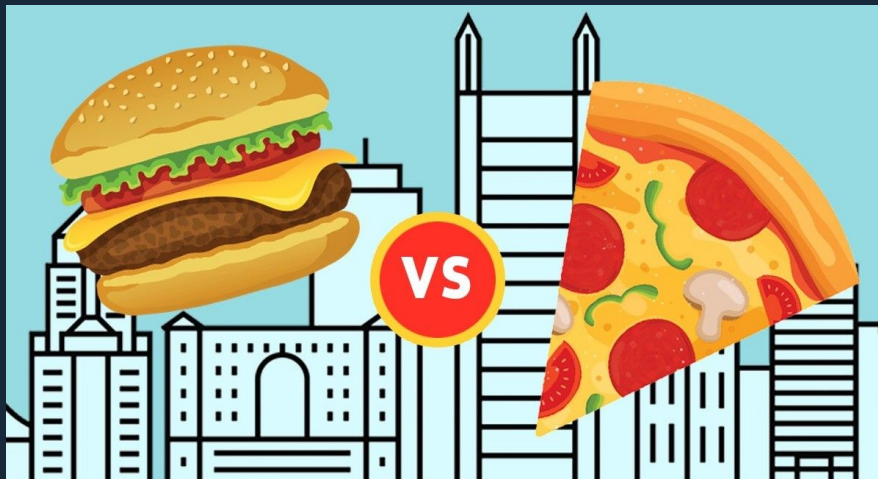
Exploit or Explore



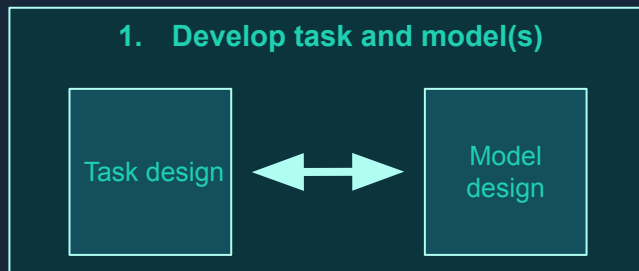
What have we learnt about choice?

Exploit **versus** Explore:

- Discover “what works” by alternating between exploration and exploitation
- In uncertain environments, more exploration **could** be useful



Summary



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