

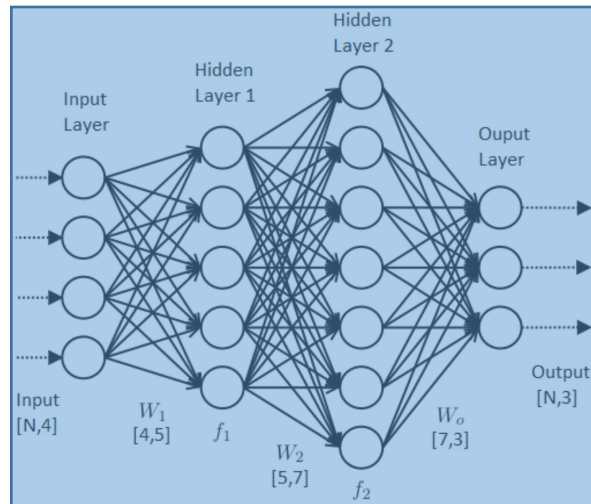
CSCSE 636 Neural Networks (Deep Learning)

Lecture 13: Deep Reinforcement Learning (continued)

Anxiao (Andrew) Jiang

Policy-based Approach

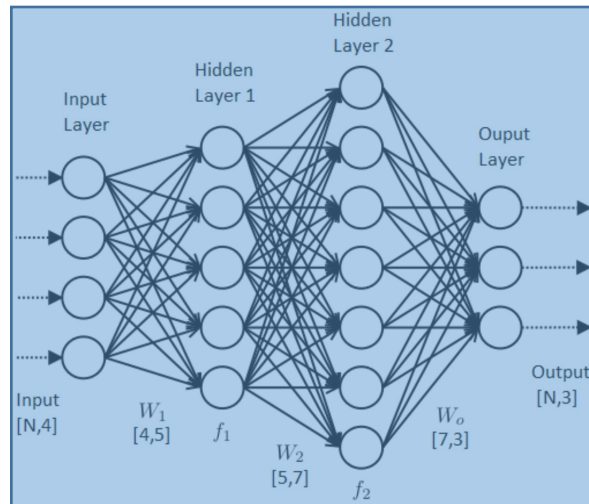
$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$



Actor



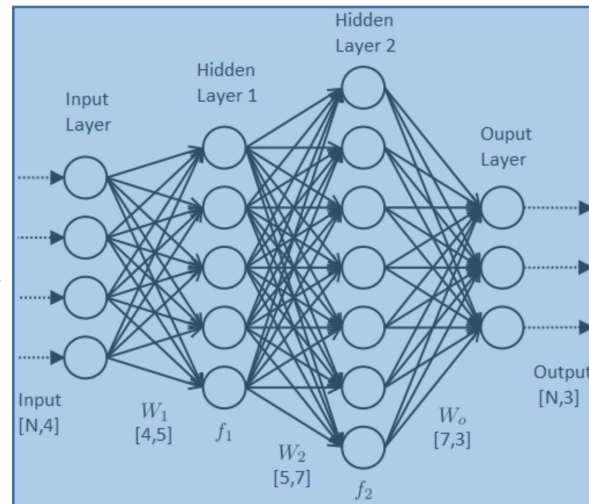
hurt leg



Actor



hurt leg



Actor

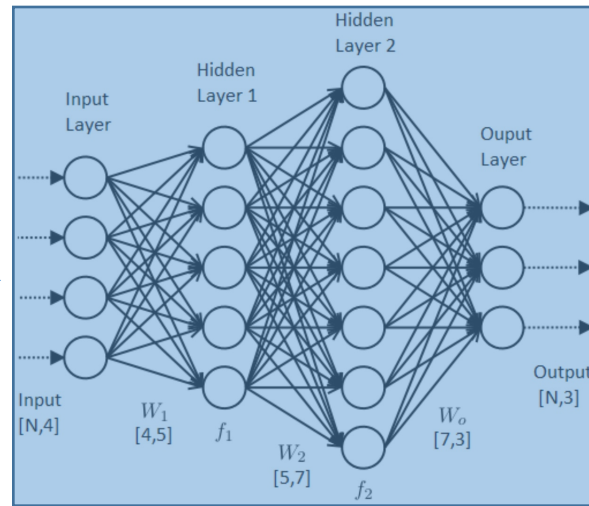
train hard

train less hard

give up



hurt leg



Actor

train hard

0.4

train less hard

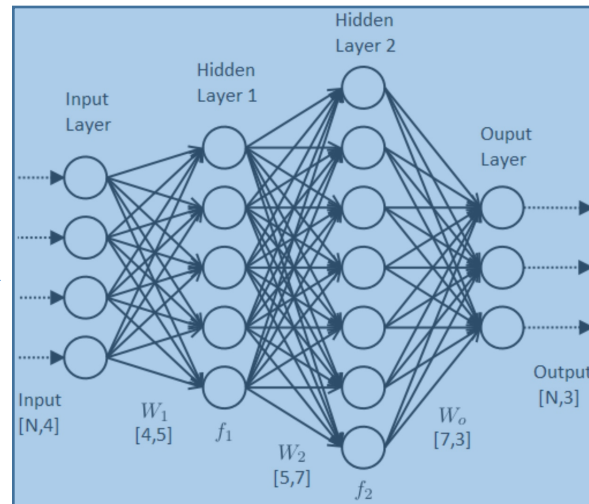
0.3

give up

0.3



hurt leg



Actor

train hard

0.4

train less hard

0.3

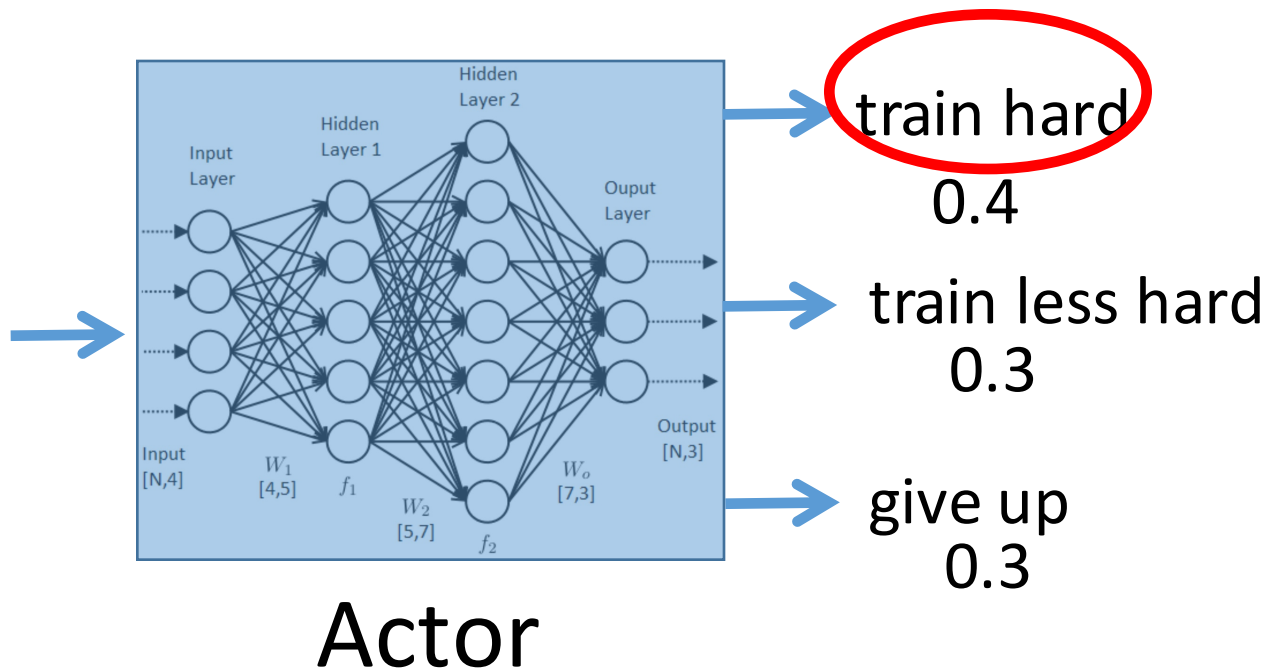
give up

0.3

Two years later ...

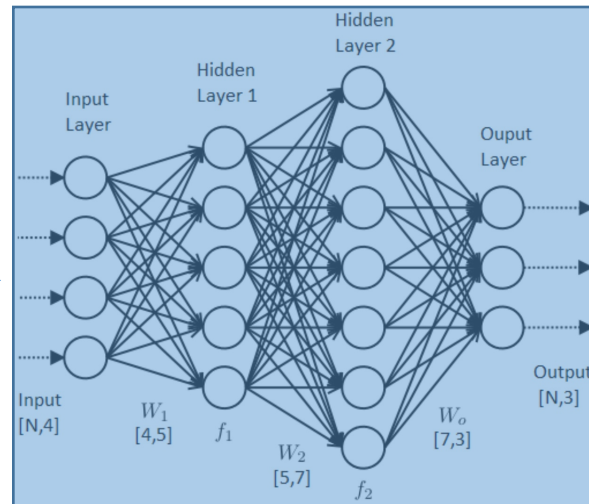


Now train the actor (neural network) ...





hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

Target
output

1

0

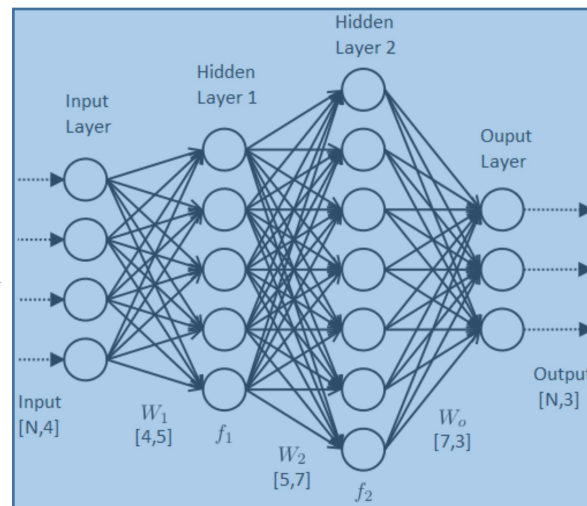
0

Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

1

0

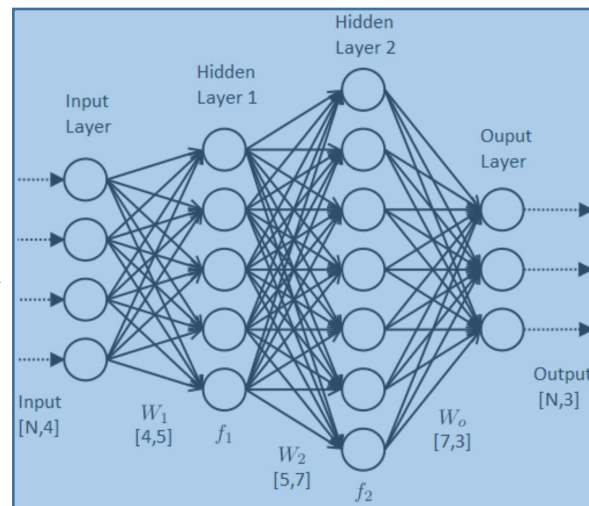
0

Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

1

0

0

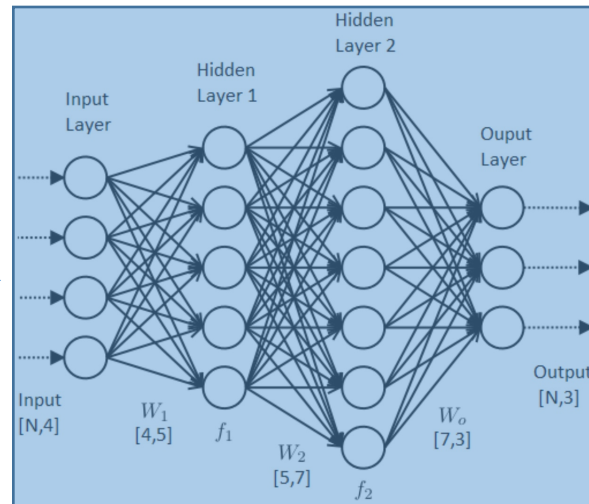
Loss = cross-entropy between the two probability vectors

Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

1

0

0

Loss = cross-entropy between the two probability vectors X

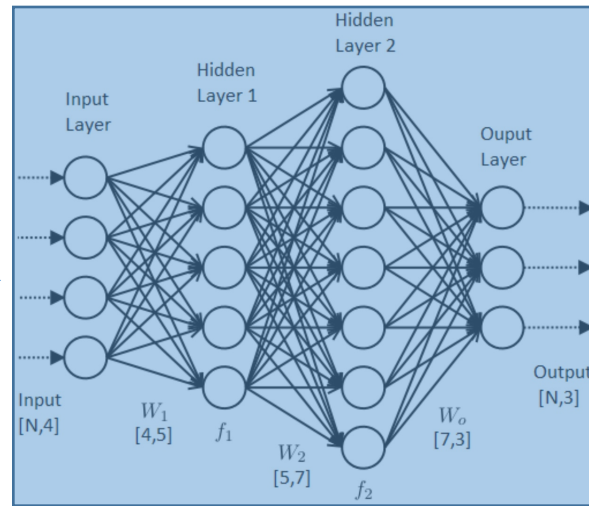


Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4 → 0.9

train less hard

0.3 → 0.06

give up

0.3 → 0.04

1

0

0

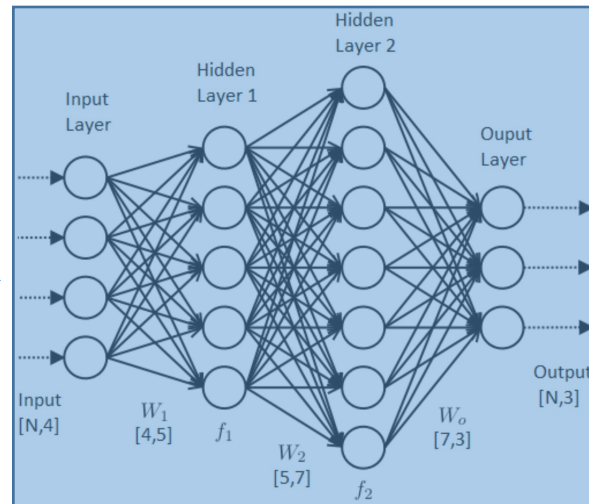
Loss = cross-entropy between the two probability vectors X



Another scenario...



hurt leg



Actor

train hard

0.4

train less hard

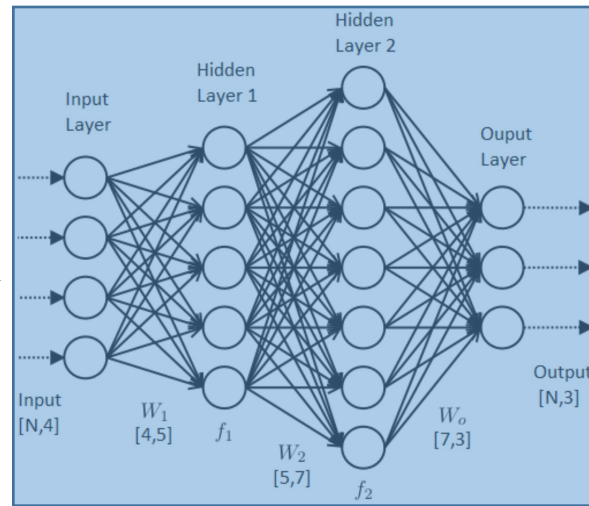
0.3

give up

0.3



hurt leg



Actor

train hard

0.4

train less hard

0.3

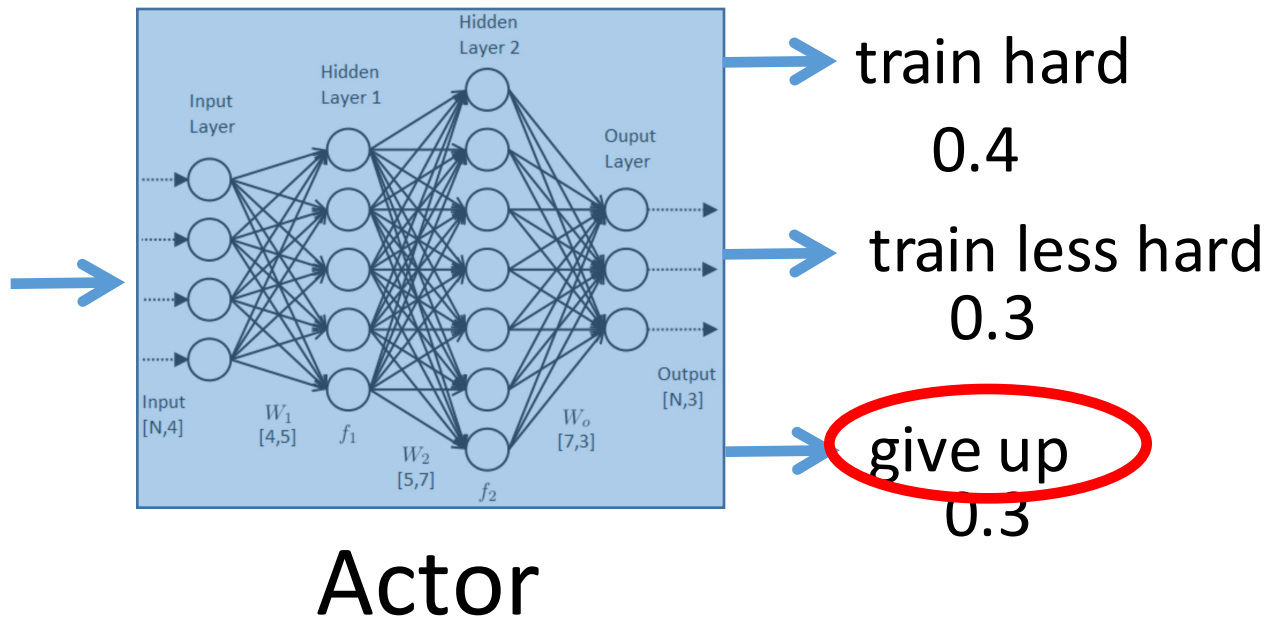
give up

0.3

Two years later ...

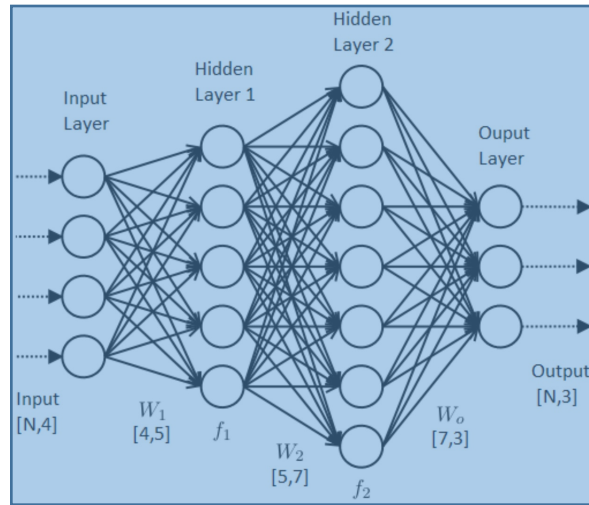


Now train the actor (neural network) ...





hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

Target
output

0

0

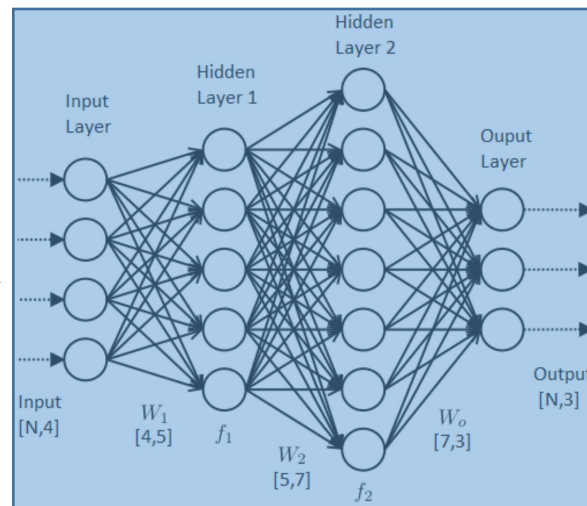
1

Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

0

0

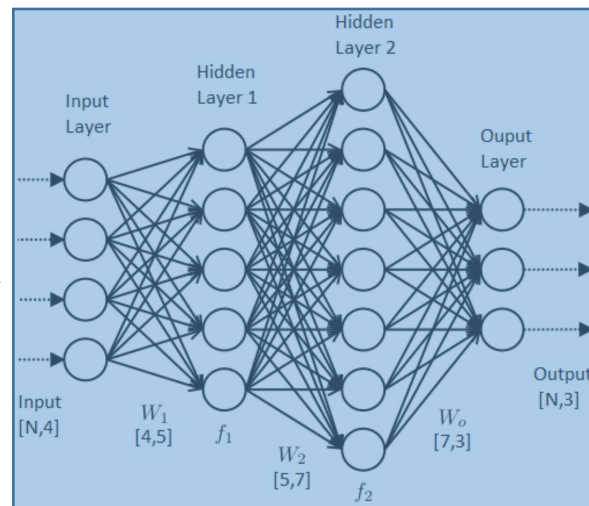
1

Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

0

0

1

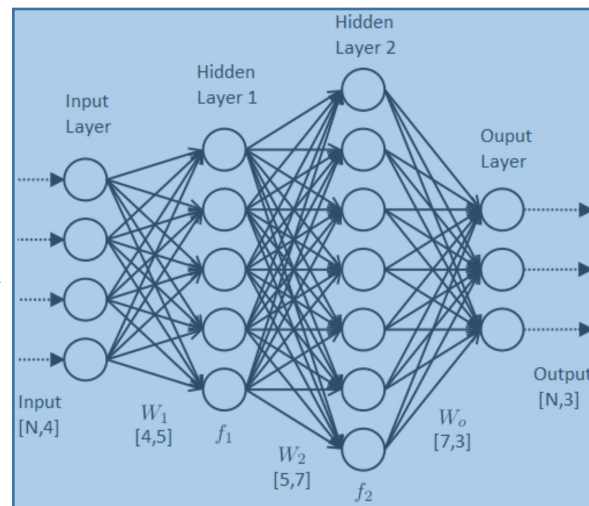
Loss = cross-entropy between the two probability vectors

Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

0

0

1

Loss = cross-entropy between the two probability vectors X

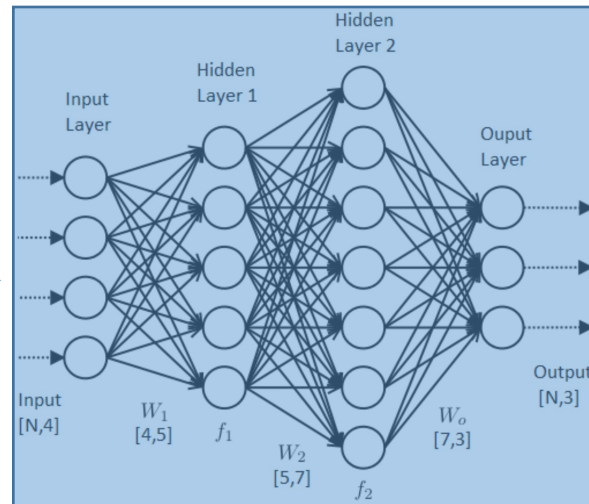


Tune weights to help the actor make the same choice with a higher probability.

Target
output



hurt leg



Actor

train hard

0.4 → 0.22

train less hard

0.3 → 0.28

give up

0.3 → 0.5

0

0

1

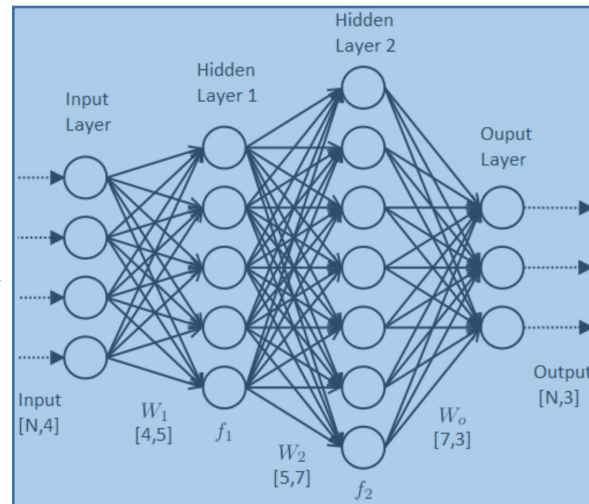
Loss = cross-entropy between the two probability vectors X



Another scenario...



hurt leg



Actor

train hard

0.4

train less hard

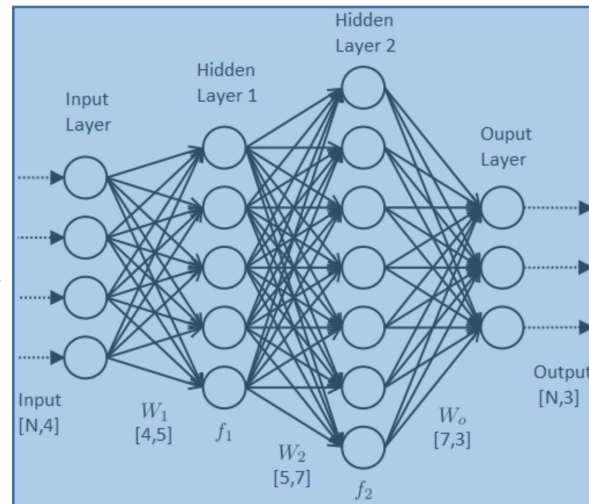
0.3

give up

0.3



hurt leg



Actor

train hard

0.4

train less hard

0.3

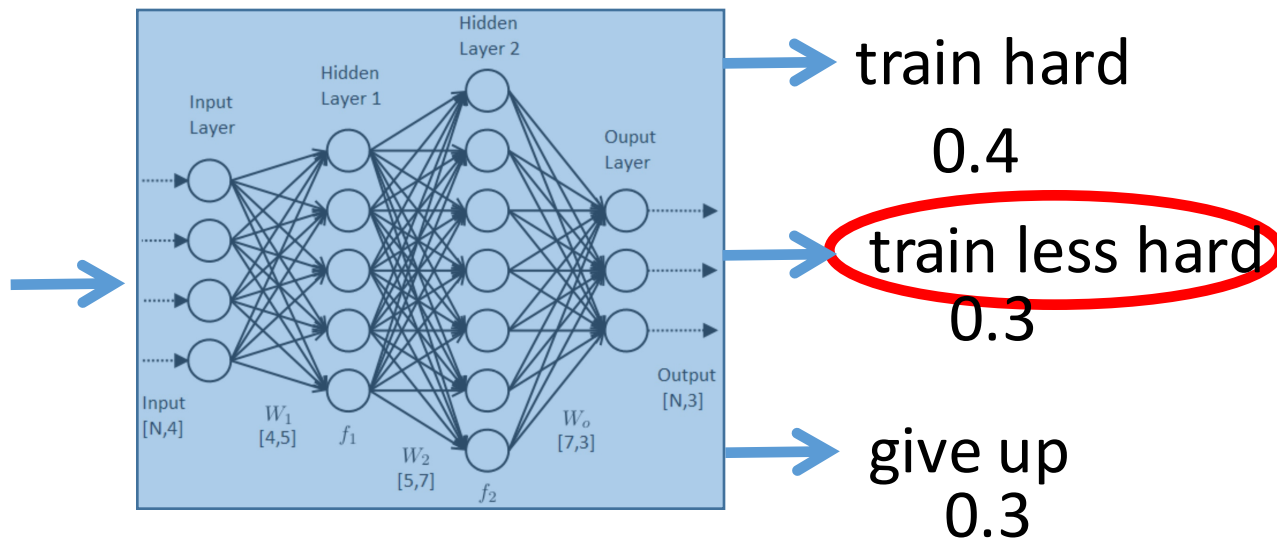
give up

0.3

Three years later ...



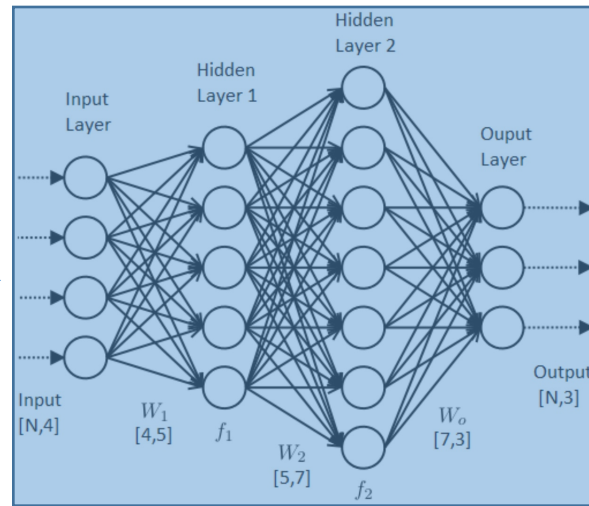
Now train the actor (neural network) ...



Actor



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

Target
output

0

1

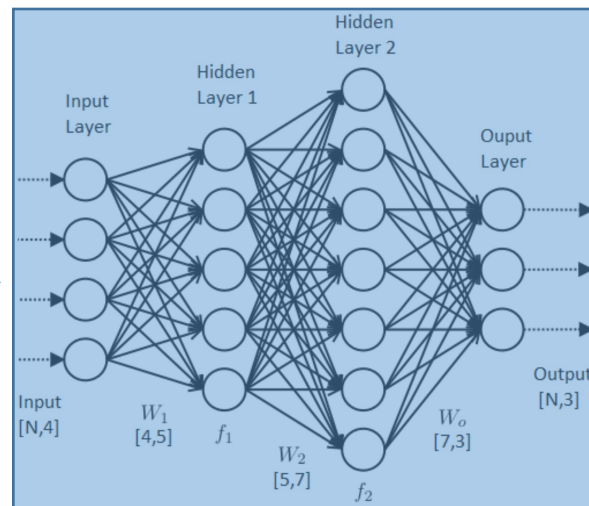
0

Tune weights to help the actor make the same choice with a lower probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up
0.3

0

1

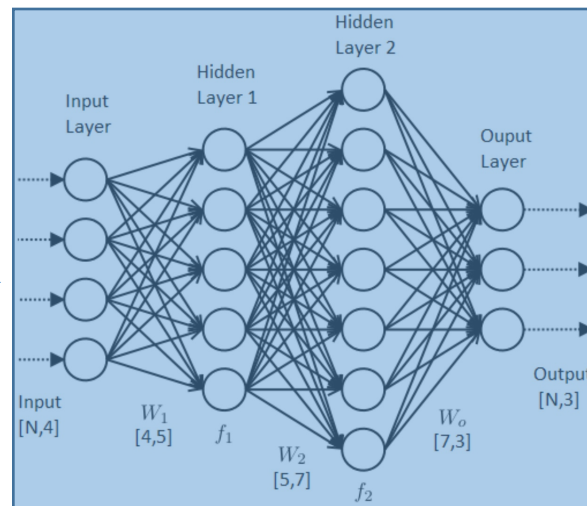
0

Tune weights to help the actor make the same choice with a lower probability.

Target
output



hurt leg



train hard

0.4

train less hard

0.3

give up
0.3

0

1

0

Actor

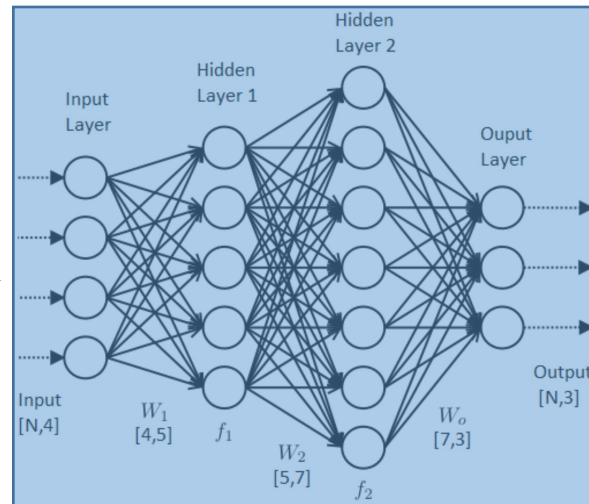
Loss = cross-entropy between the two probability vectors

Tune weights to help the actor make the same choice with a lower probability.

Target
output



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

0

1

0

Loss = cross-entropy between the two probability vectors X

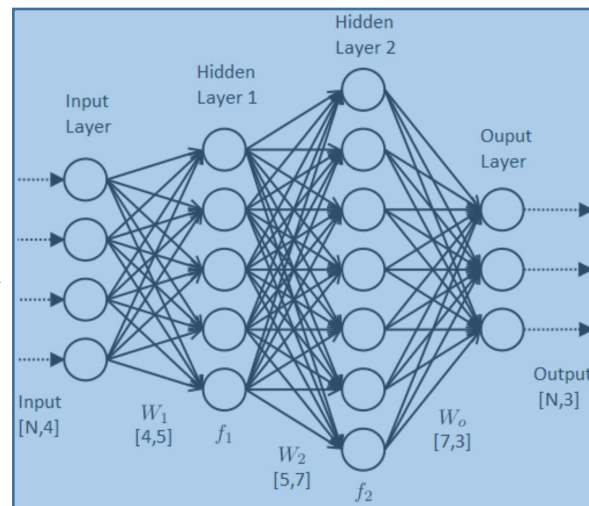


Tune weights to help the actor make the same choice with a lower probability.

Target
output



hurt leg



Actor

train hard

0.4 → 0.42

train less hard

0.3 → 0.2

give up

0.3 → 0.38

0

1

0

Loss = cross-entropy between the two probability vectors X

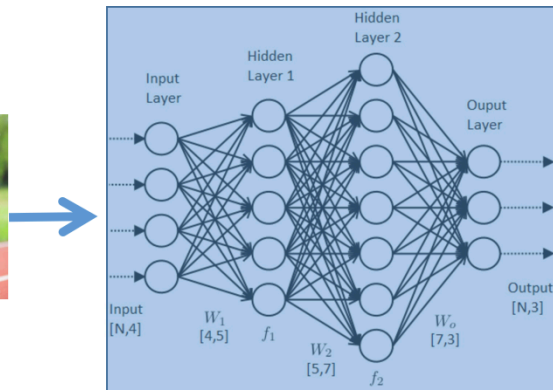


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3



$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$



negative
Cross
entropy

Target
output

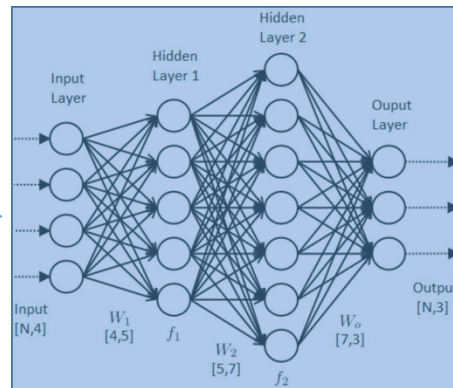
1

0

0



hurt leg



train hard

0.4

train less hard

0.3

give up

0.3

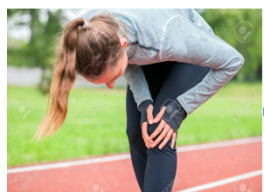
Actor

Minimize Cross Entropy:

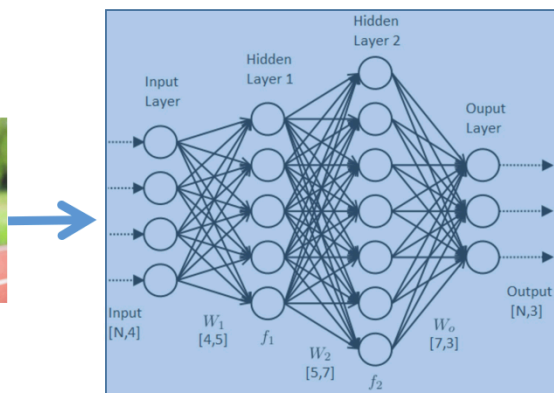
$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

negative
Cross
entropy

Target
output

1

0

0



$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$

↑
negative
Cross
entropy

Target
output

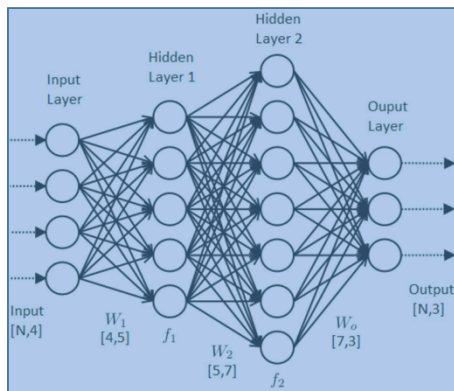
1

0

0



hurt leg



train hard

0.4

train less hard

0.3

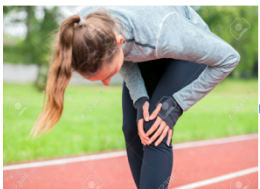
give up

0.3

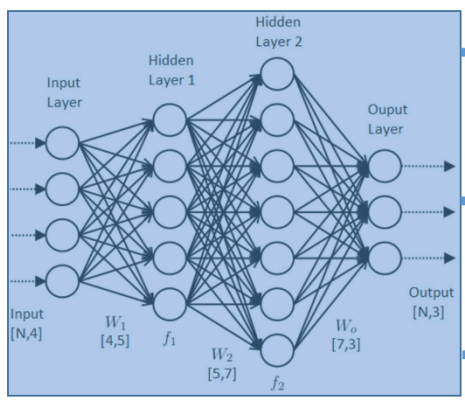
Actor

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

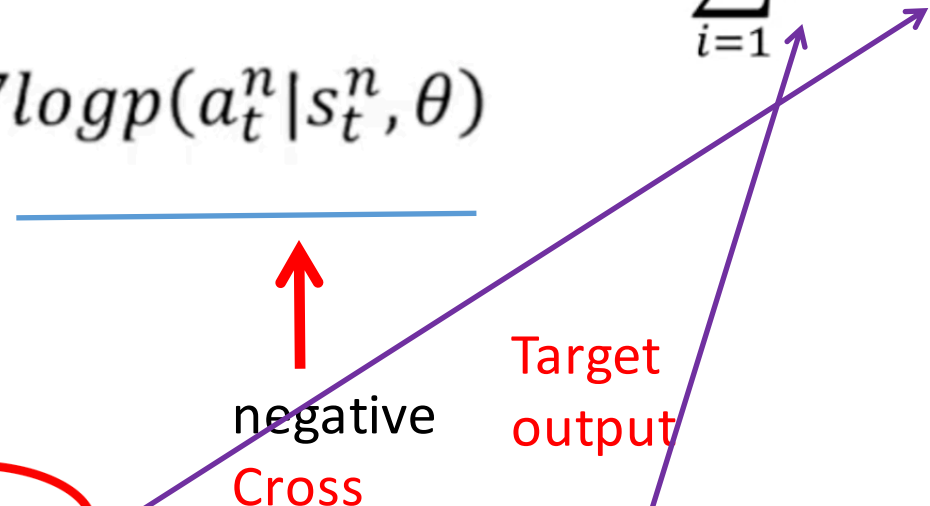
negative
Cross
entropy

Target
output

1

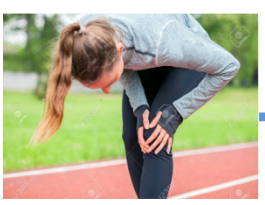
0

0

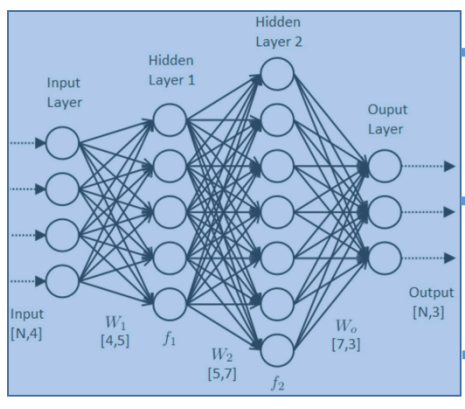


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$



hurt leg



Actor

train hard
0.4

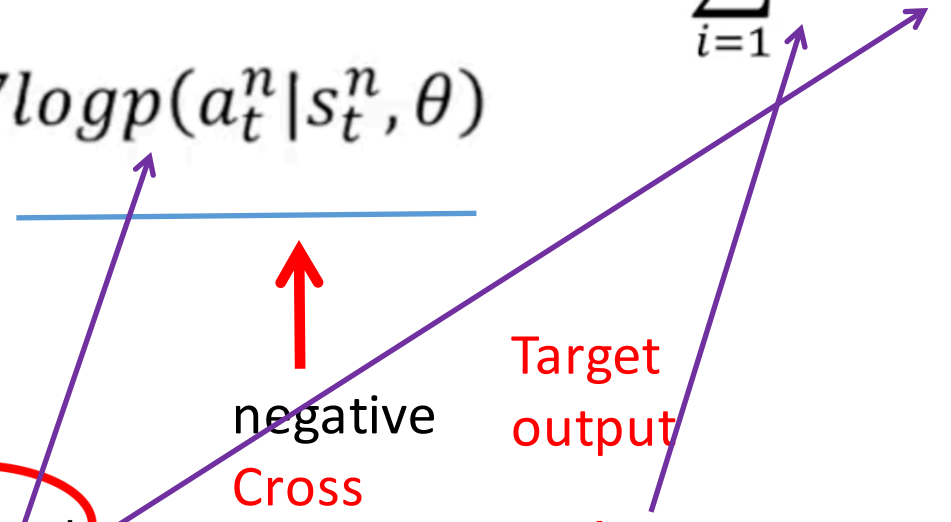
train less hard
0.3

give up
0.3

negative
Cross
entropy

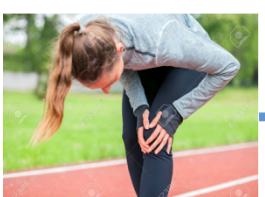
Target
output

1
0
0

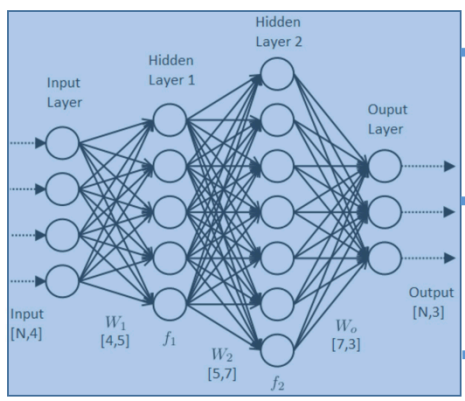


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

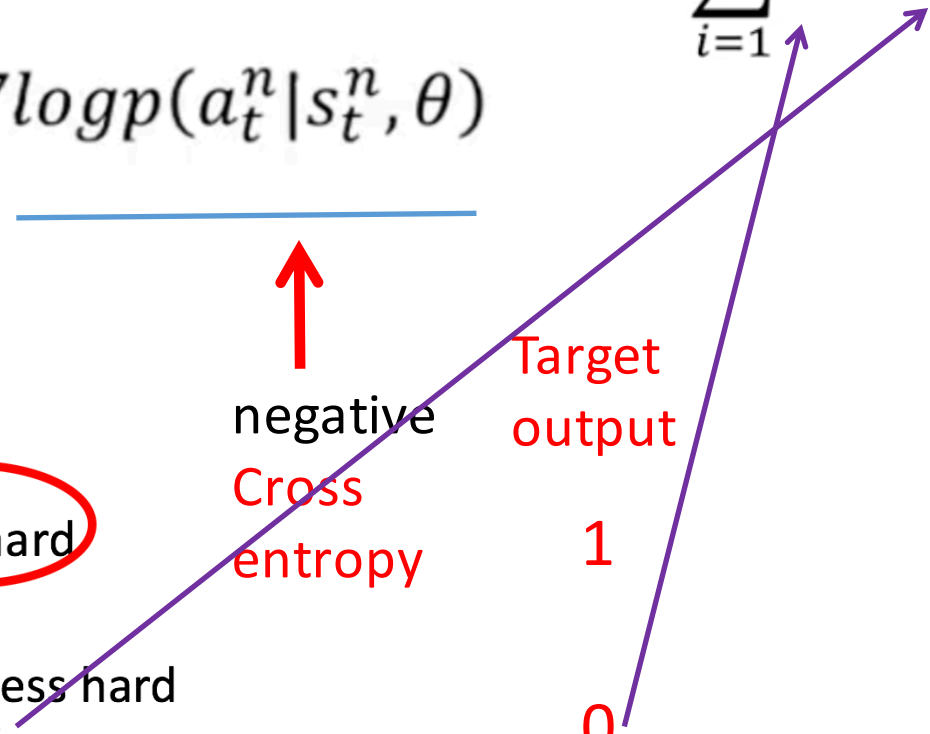
negative
Cross
entropy

Target
output

1

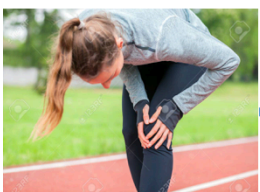
0

0

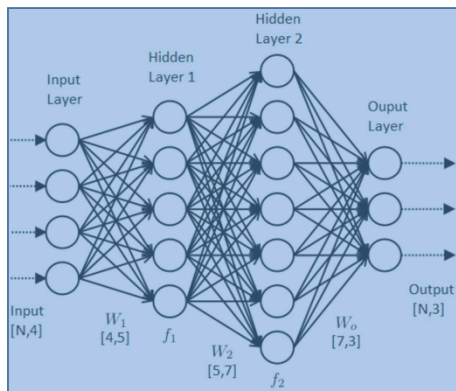


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

$$-\sum_{i=1}^3 \hat{y}_i \log y_i$$



hurt leg



Actor

train hard

0.4

train less hard

0.3

give up

0.3

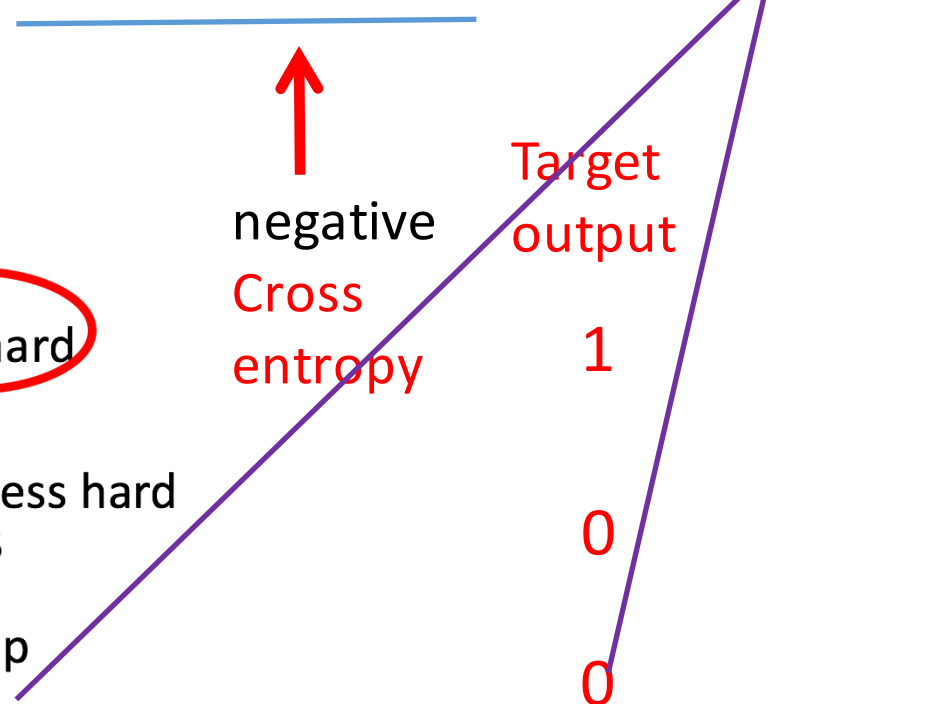
negative
Cross
entropy

Target
output

1

0

0

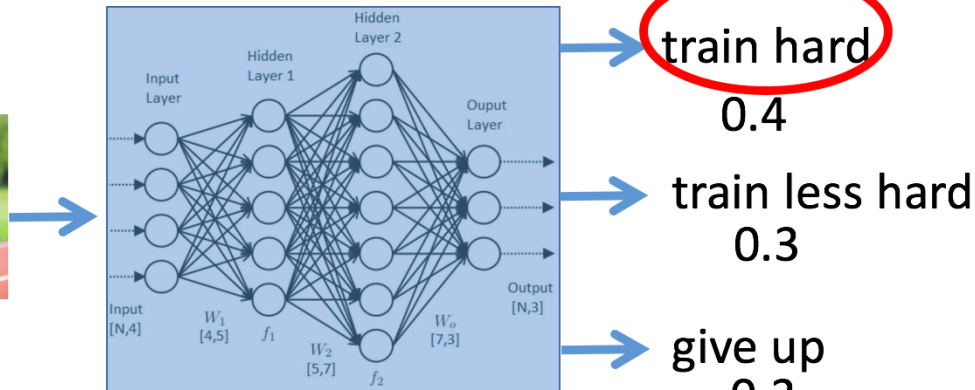


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \underbrace{(R(\tau^n) - b)}_{\text{Constant weight for this input-output pair}} \nabla \log p(a_t^n | s_t^n, \theta)$$

Constant weight for
this input-output pair



hurt leg



Actor

Target
output

1

0

0

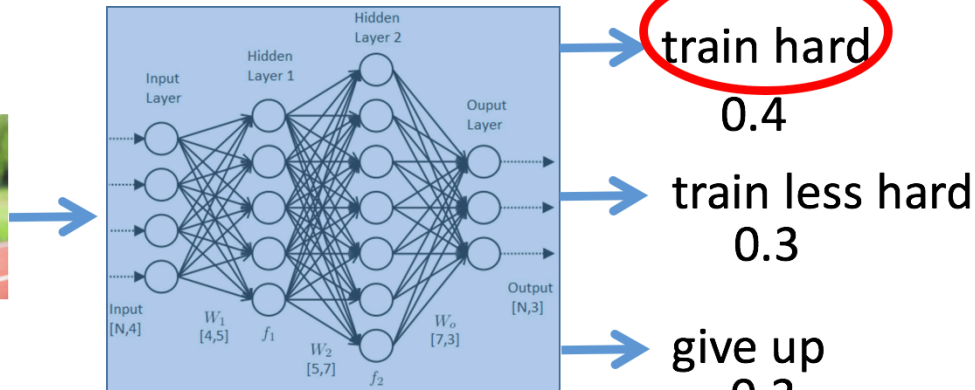
$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$



Sum over all the input-output pairs



hurt leg



Actor

Target
output

1

0

0

train hard

0.4

train less hard

0.3

give up

0.3

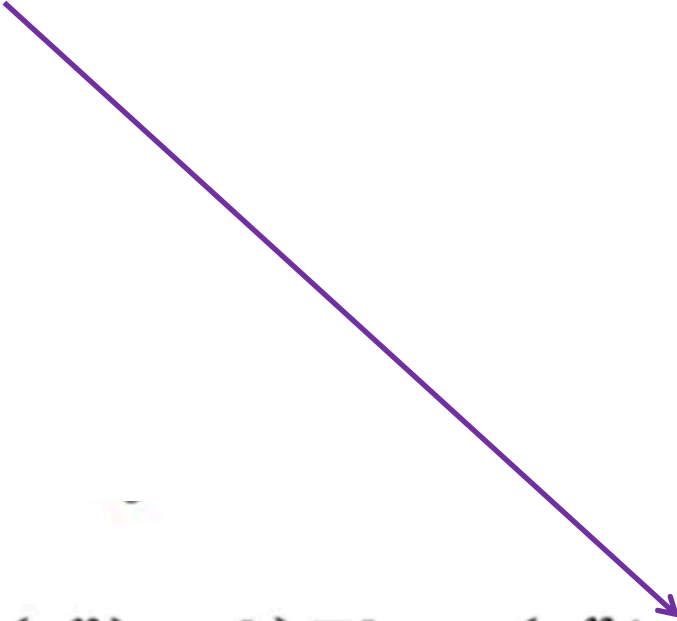
1. Use the current neural network (actor) to play the game, to get data from many episodes.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

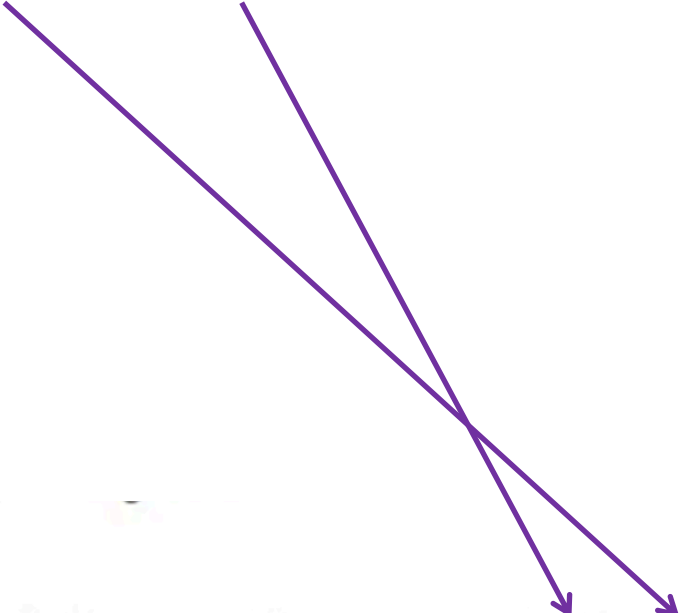
1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

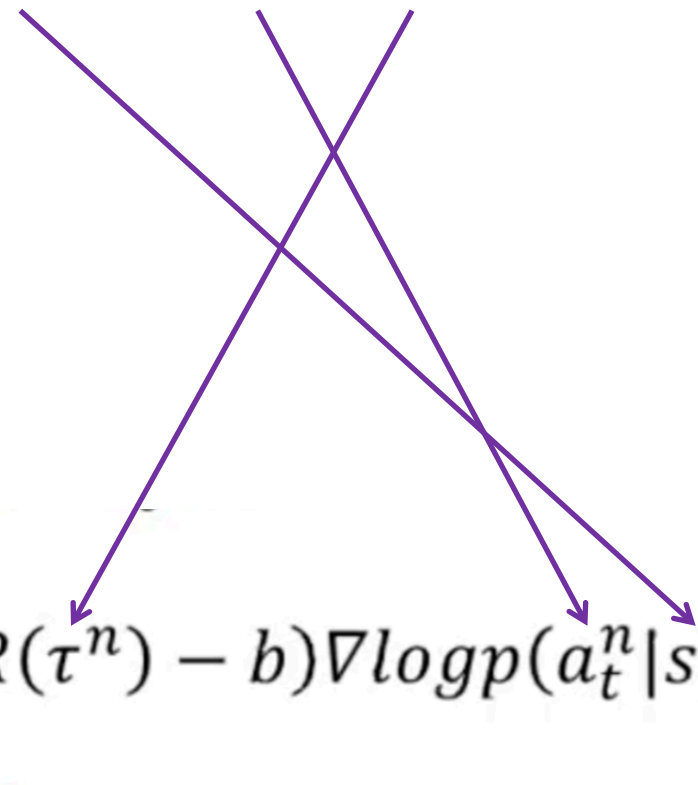
1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$


1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).

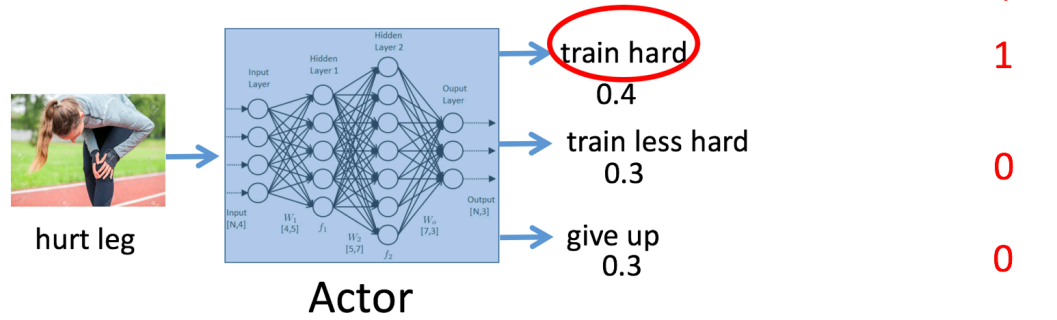


The diagram shows three purple arrows originating from the text above. One arrow points from the word 'reward' in the second step to the term $R(\tau^n)$ in the equation. Another arrow points from the word 'action' in the second step to the term a_t^n in the equation. A third arrow points from the word 'observation' in the second step to the term s_t^n in the equation.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).
3. Turn the reinforcement learning problem into a classification problem. Train the network.

Sum over all the input-output pairs



$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{\tau^n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).
3. Turn the reinforcement learning problem into a classification problem. Train the network.
4. Use the new network to collect more data. Use the new data to train the new network.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).
3. Turn the reinforcement learning problem into a classification problem. Train the network.
4. Use the new network to collect more data. Use the new data to train the new network.
5. Repeat the above steps, until the network's performance converges.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

The following lecture is based on the interesting lecture of Prof. Hung-yi Lee “Deep Reinforcement Learning”
https://www.youtube.com/watch?v=W8XF3ME8G2I&list=PLJV_eI3uVTsPy9oCRY30oBPNLCo89yu49&index=33

Policy Gradient

Policy Gradient

Given actor parameter θ

Policy Gradient

Given actor parameter θ

$$\tau^1: \begin{array}{l} (s_1^1, a_1^1) \\ (s_2^1, a_2^1) \\ \vdots \end{array}$$

Policy Gradient

Given actor parameter θ

$$\tau^1: \begin{array}{ll} (s_1^1, a_1^1) & R(\tau^1) \\ (s_2^1, a_2^1) & R(\tau^1) \\ \vdots & \vdots \end{array}$$

Policy Gradient

Given actor parameter θ

$$\tau^1: (s_1^1, a_1^1) \quad R(\tau^1)$$

$$(s_2^1, a_2^1) \quad R(\tau^1)$$

$$\vdots$$
$$\vdots$$

$$\tau^2: (s_1^2, a_1^2) \quad R(\tau^2)$$

$$(s_2^2, a_2^2) \quad R(\tau^2)$$

$$\vdots$$
$$\vdots$$

Policy Gradient

Given actor parameter θ

$$\begin{array}{l} \tau^1: (s_1^1, a_1^1) \quad R(\tau^1) \\ \quad (s_2^1, a_2^1) \quad R(\tau^1) \\ \quad \vdots \quad \quad \quad \vdots \\ \tau^2: (s_1^2, a_1^2) \quad R(\tau^2) \\ \quad (s_2^2, a_2^2) \quad R(\tau^2) \\ \quad \vdots \quad \quad \quad \vdots \end{array}$$

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

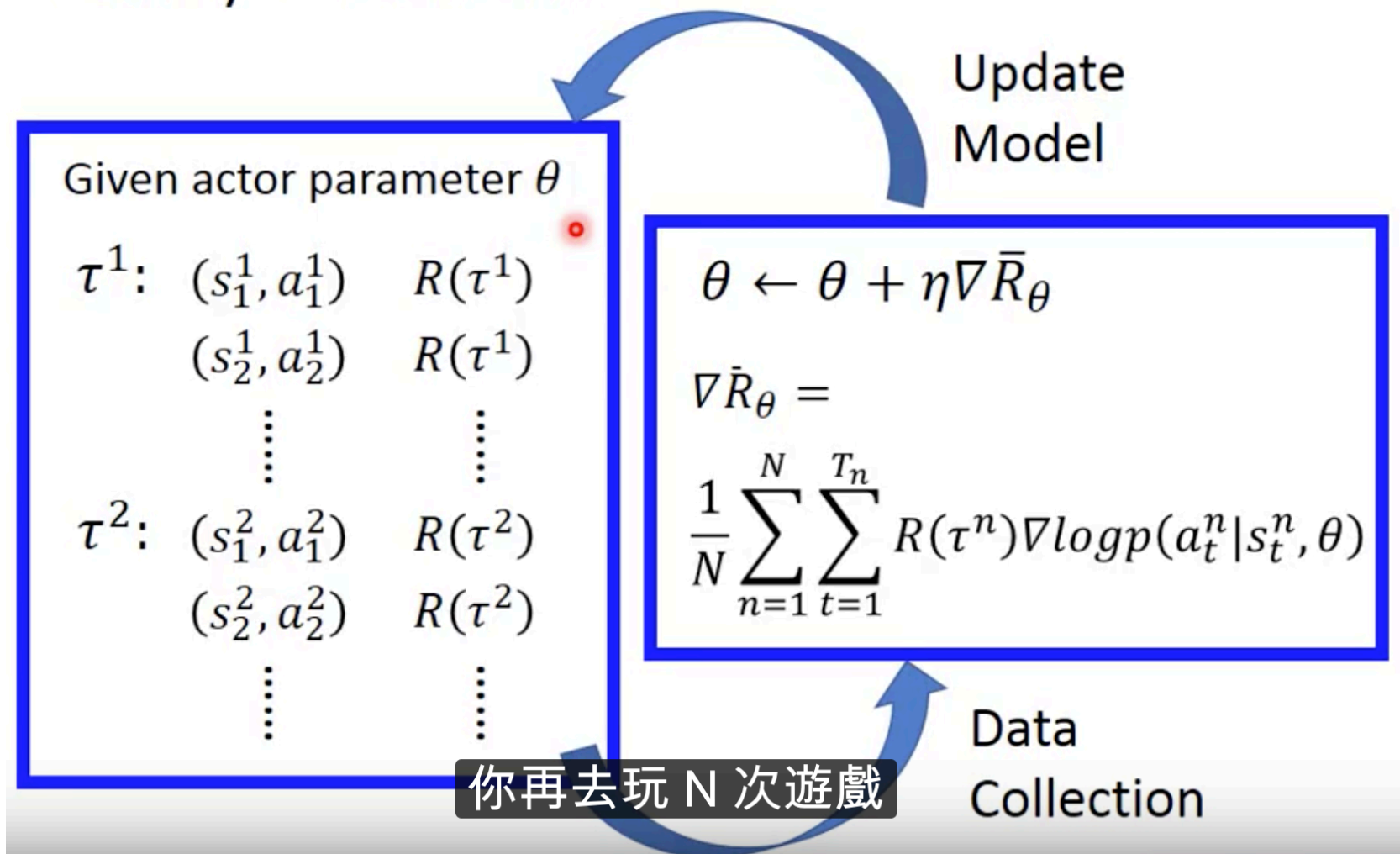
$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta)$$

推出來的這個式子

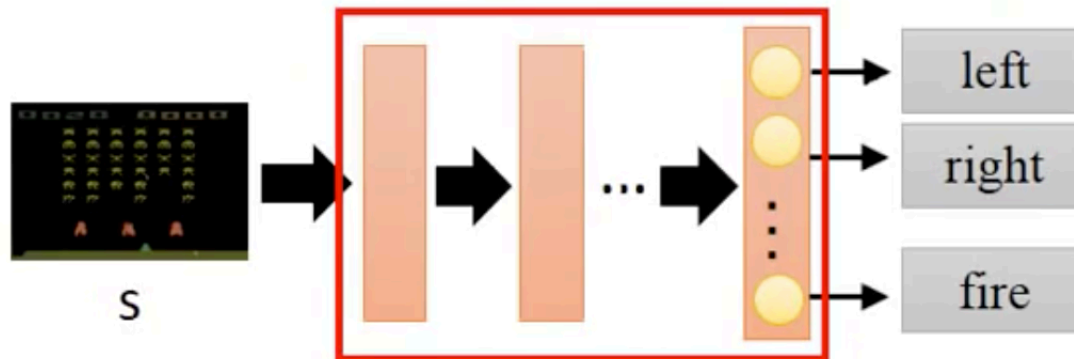
Data
Collection

Policy Gradient



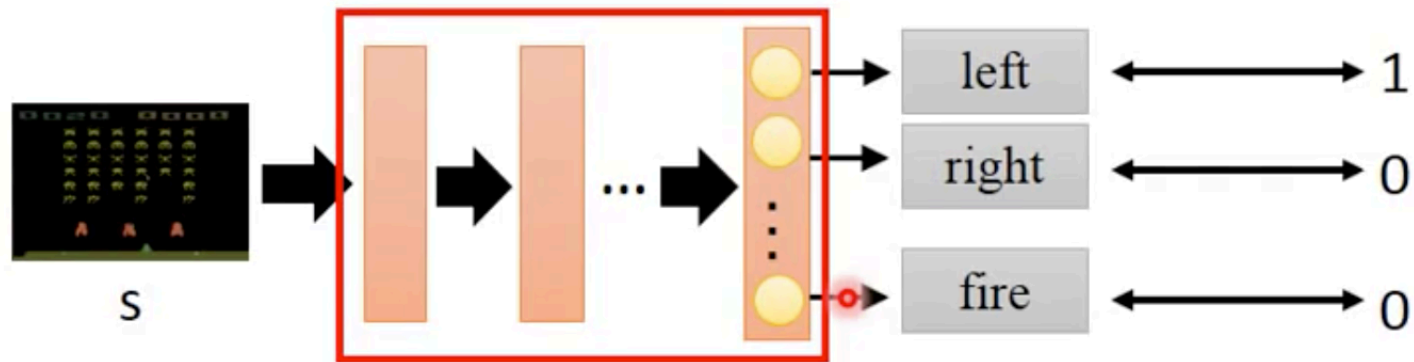
Policy Gradient

Considered as
Classification Problem



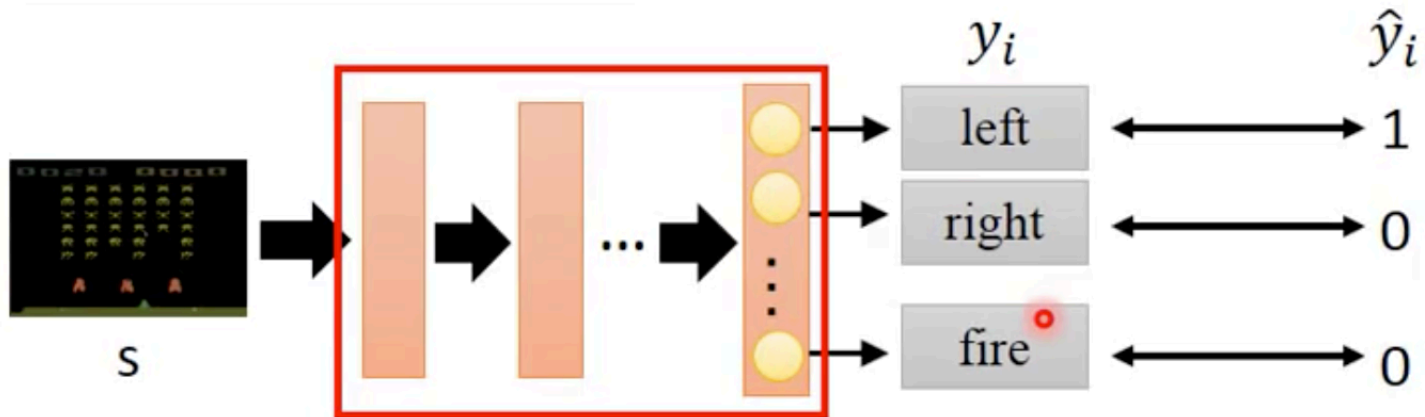
Policy Gradient

Considered as
Classification Problem



Policy Gradient

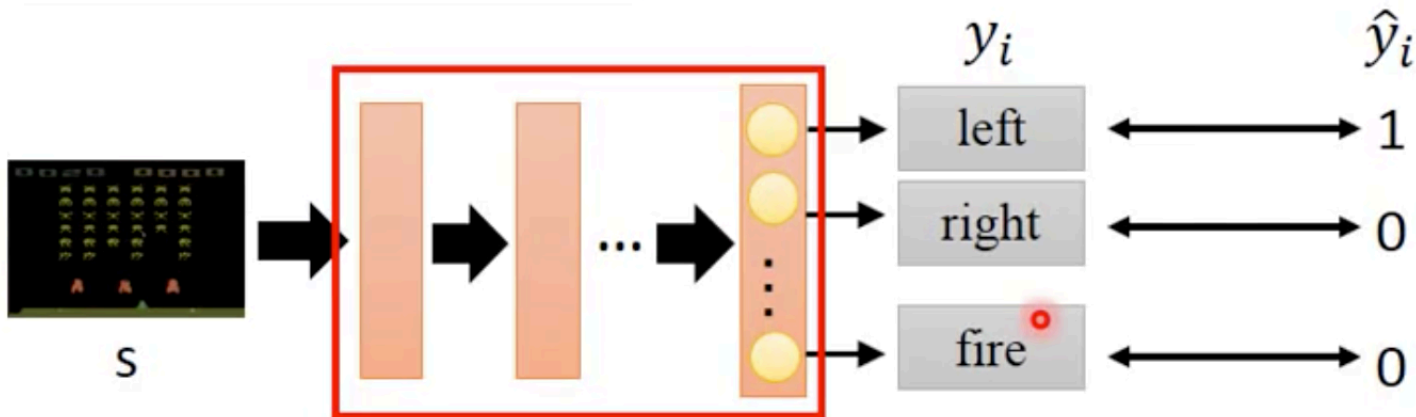
Considered as
Classification Problem



Policy Gradient

Considered as
Classification Problem

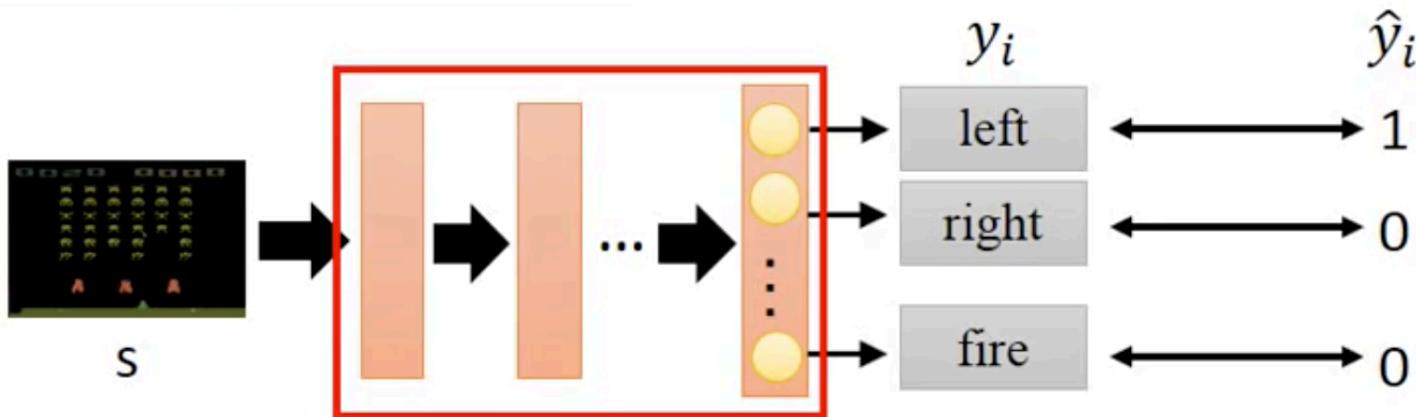
$$\text{Minimize: } - \sum_{i=1}^3 \hat{y}_i \log y_i$$



Policy Gradient

Considered as
Classification Problem

Minimize: $-\sum_{i=1}^3 \hat{y}_i \log y_i$

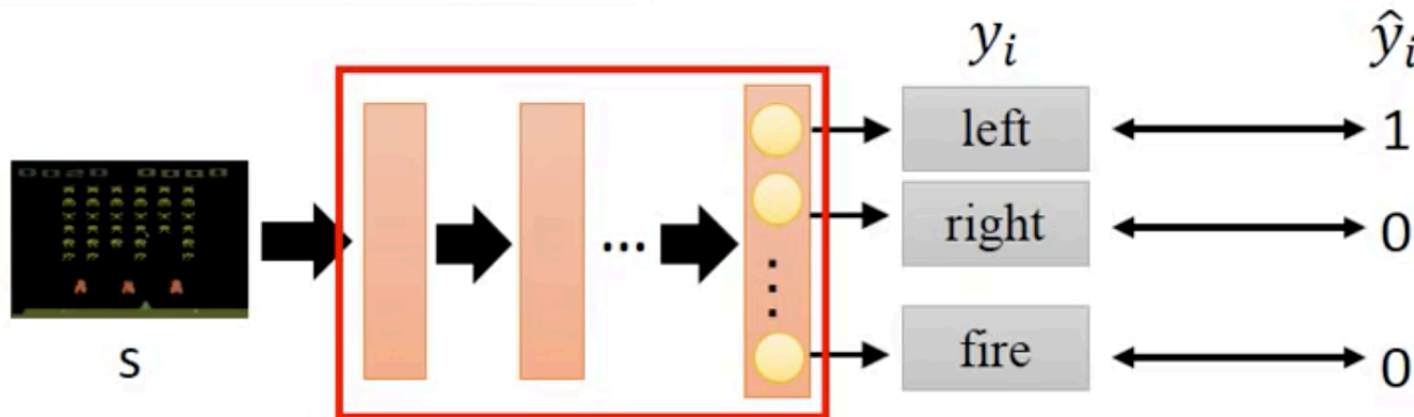


Maximize: $\log y_i =$

Policy Gradient

Considered as
Classification Problem

$$\text{Minimize: } - \sum_{i=1}^3 \hat{y}_i \log y_i$$

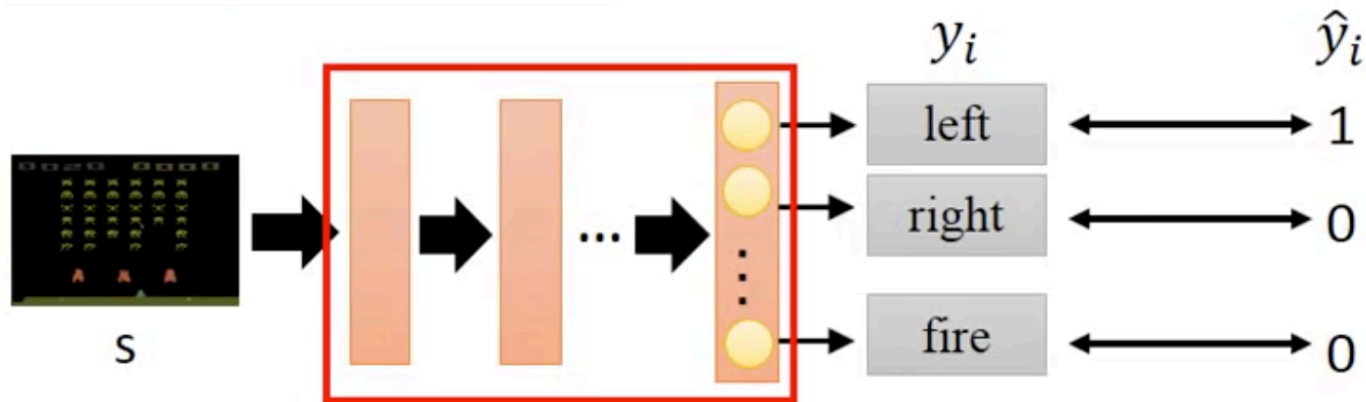


$$\text{Maximize: } \log y_i = \log P(\text{"left"}|s)$$

Policy Gradient

Considered as
Classification Problem

$$\text{Minimize: } - \sum_{i=1}^3 \hat{y}_i \log y_i$$



$$\text{Maximize: } \log y_i = \log P(\text{"left"}|s)$$

$$\theta \leftarrow \theta + \eta \nabla \log P(\text{"left"}|s)$$

Policy Gradient

Given actor parameter θ

τ^1 :	(s_1^1, a_1^1)	$R(\tau^1)$
	(s_2^1, a_2^1)	$R(\tau^1)$
	\vdots	\vdots
τ^2 :	(s_1^2, a_1^2)	$R(\tau^2)$
	(s_2^2, a_2^2)	$R(\tau^2)$
	\vdots	\vdots

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \boxed{} \nabla \log p(a_t^n | s_t^n, \theta)$$

Policy Gradient

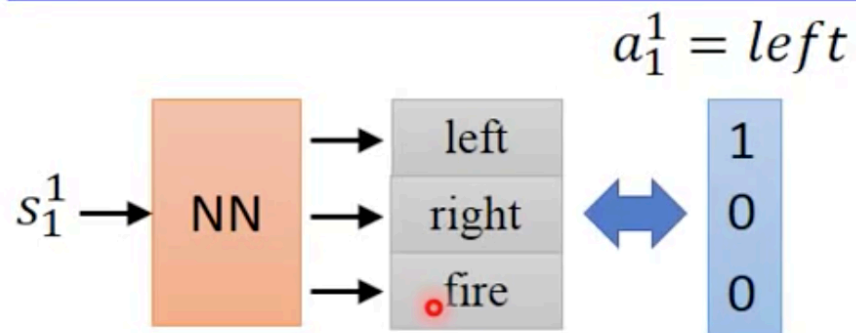
Given actor parameter θ

τ^1 : (s_1^1, a_1^1) $R(\tau^1)$
 (s_2^1, a_2^1) $R(\tau^1)$
 \vdots \vdots
 τ^2 : (s_1^2, a_1^2) $R(\tau^2)$
 (s_2^2, a_2^2) $R(\tau^2)$
 \vdots \vdots

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \blacksquare \nabla \log p(a_t^n | s_t^n, \theta)$$



Policy Gradient

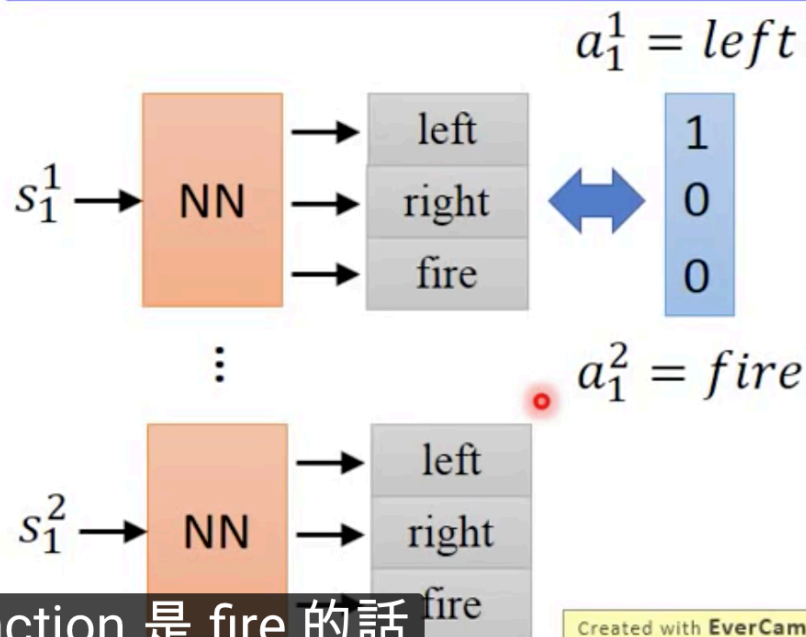
Given actor parameter θ

$\tau^1: (s_1^1, a_1^1) \quad R(\tau^1)$
 $(s_2^1, a_2^1) \quad R(\tau^1)$
 $\vdots \quad \vdots$
 $\tau^2: (s_1^2, a_1^2) \quad R(\tau^2)$
 $(s_2^2, a_2^2) \quad R(\tau^2)$
 $\vdots \quad \vdots$

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \blacksquare \nabla \log p(a_t^n | s_t^n, \theta)$$



我們採取的 action 是 fire 的話

Policy Gradient

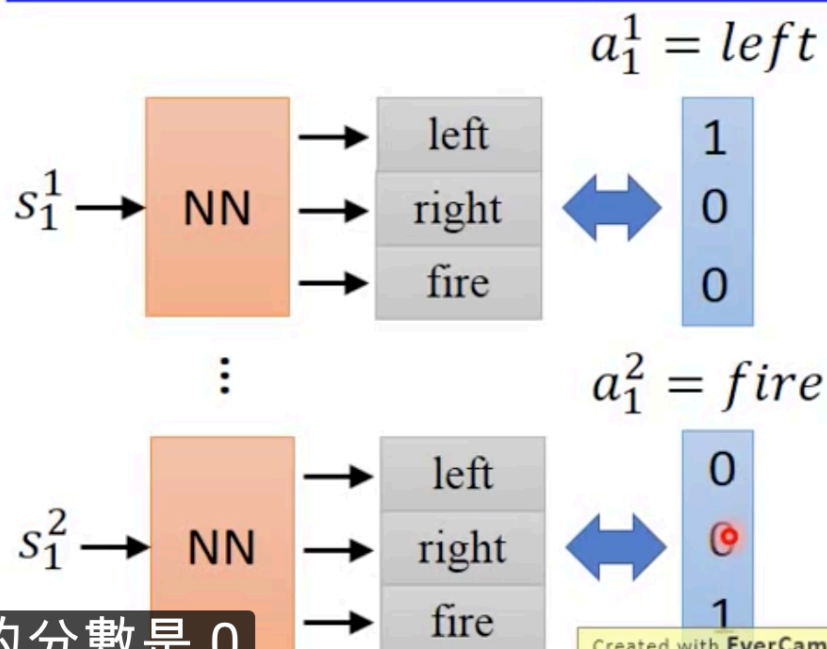
Given actor parameter θ

$\tau^1: (s_1^1, a_1^1) \quad R(\tau^1)$
 $(s_2^1, a_2^1) \quad R(\tau^1)$
 $\vdots \quad \vdots$
 $\tau^2: (s_1^2, a_1^2) \quad R(\tau^2)$
 $(s_2^2, a_2^2) \quad R(\tau^2)$
 $\vdots \quad \vdots$

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \square \nabla \log p(a_t^n | s_t^n, \theta)$$



Policy Gradient

Given actor parameter θ

$$\begin{array}{l} \tau^1: \quad (s_1^1, a_1^1) \quad R(\tau^1) \\ \quad \quad (s_2^1, a_2^1) \quad R(\tau^1) \\ \quad \quad \quad \vdots \quad \quad \quad \vdots \\ \tau^2: \quad (s_1^2, a_1^2) \quad R(\tau^2) \\ \quad \quad (s_2^2, a_2^2) \quad R(\tau^2) \\ \quad \quad \quad \vdots \quad \quad \quad \vdots \end{array}$$

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta)$$

Policy Gradient

Given actor parameter θ

$\tau^1:$	(s_1^1, a_1^1)	$R(\tau^1)$
	(s_2^1, a_2^1)	$R(\tau^1)$
	\vdots	\vdots
$\tau^2:$	(s_1^2, a_1^2)	$R(\tau^2)$
	(s_2^2, a_2^2)	$R(\tau^2)$
	\vdots	\vdots

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta)$$

Each training data is weighted by $R(\tau^n)$

Policy Gradient

Given actor parameter θ

τ^1 :	(s_1^1, a_1^1)	$R(\tau^1)$	2
	(s_2^1, a_2^1)	$R(\tau^1)$	2
	\vdots	\vdots	
τ^2 :	(s_1^2, a_1^2)	$R(\tau^2)$	1
	(s_2^2, a_2^2)	$R(\tau^2)$	1
	\vdots	\vdots	

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta)$$

Each training data is weighted by $R(\tau^n)$

Policy Gradient

Given actor parameter θ

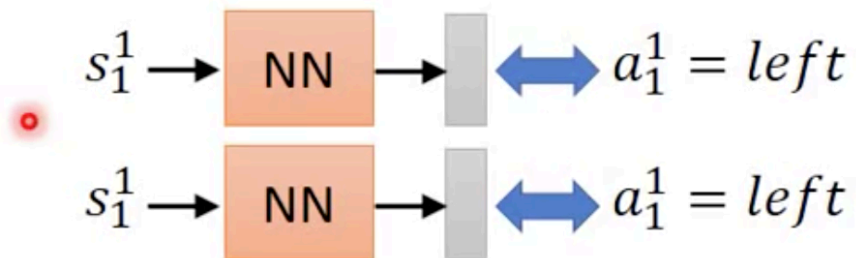
τ^1 :	(s_1^1, a_1^1)	$R(\tau^1)$	2
	(s_2^1, a_2^1)	$R(\tau^1)$	2
	\vdots	\vdots	
τ^2 :	(s_1^2, a_1^2)	$R(\tau^2)$	1
	(s_2^2, a_2^2)	$R(\tau^2)$	1
	\vdots	\vdots	

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta)$$

Each training data is weighted by $R(\tau^n)$



Policy Gradient

Given actor parameter θ

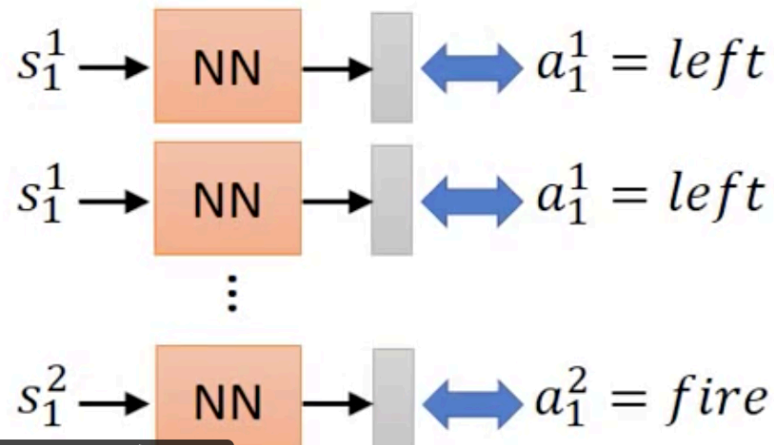
τ^1 :	(s_1^1, a_1^1)	$R(\tau^1)$	2
	(s_2^1, a_2^1)	$R(\tau^1)$	2
	\vdots	\vdots	
τ^2 :	(s_1^2, a_1^2)	$R(\tau^2)$	1
	(s_2^2, a_2^2)	$R(\tau^2)$	1
	\vdots	\vdots	

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta)$$

Each training data is weighted by $R(\tau^n)$



1. Use the current neural network (actor) to play the game, to get data from many episodes.
2. As a result, we get many triplets (observation, action, reward).
3. Turn the reinforcement learning problem into a classification problem. Train the network.
4. Use the new network to collect more data. Use the new data to train the new network.
5. Repeat the above steps, until the network's performance converges.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$