

---

# **Off-policy control with approximation**

---

**Herke van Hoof**

# Today's chapter

---

This week will be a tricky subject - don't stress out if you don't understand everything. From next lecture, we'll start on a new topic, starting a bit easier.

That said, do make use of the resources (exercises, book, slides, tutorial sessions) to review the material as necessary.

---

# Control with approximation

---

Last lecture: Semi-gradient SARSA

Straightforward analogue of semi-gradient TD:

Interleave using  $\epsilon$ -greedy policy with weight update

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha (R_t + \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}_t) - \hat{q}(S_t, A_t, \mathbf{w}_t)) \nabla \hat{q}(S_t, A_t, \mathbf{w}_t)$$

# Off-policy prediction with approximation

---

That was quite straightforward so far!

How about off-policy?

Remember importance weights  $\rho_t = \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$

Sample ‘too rare’: high weight

Sample ‘too common’: low weight

# Off-policy prediction with approximation

---

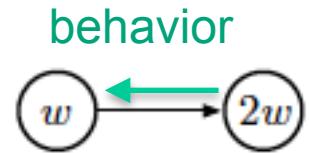
Use importance weight in semi-gradient TD(0)

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha \rho_t \underbrace{(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t))}_{\delta_t} \nabla \hat{v}(S_t, \mathbf{w}_t)$$

Let's try it!

# Off-policy prediction with approximation

Initial  $w$  is 10



On-policy,  $\gamma=1$

From left to right:

$$\begin{aligned} w_{t+1} &= w_t + \alpha(2w_t - w_t)\nabla\hat{v} \\ &= w_t + \alpha(2 - 1)w_t \end{aligned}$$

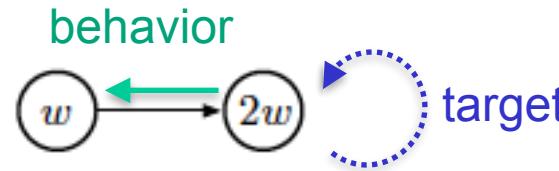
Right to left:

$$\begin{aligned} w_{t+1} &= w_t + \alpha(w_t - 2w_t)\nabla\hat{v} \\ &= w_t + \alpha(1 - 2)w_t \cdot 2 \end{aligned}$$

Mean TDE tends to shrink!

# Off-policy prediction with approximation

Initial  $w$  is 10



On-policy,  $\gamma=1$

From left to right:

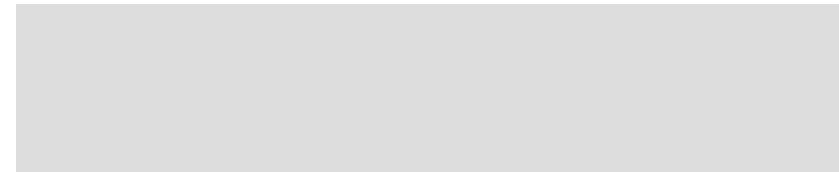
$$\begin{aligned} w_{t+1} &= w_t + \alpha(2w_t - w_t)\nabla\hat{v} \\ &= w_t + \alpha(2 - 1)w_t \end{aligned}$$

Right to left:

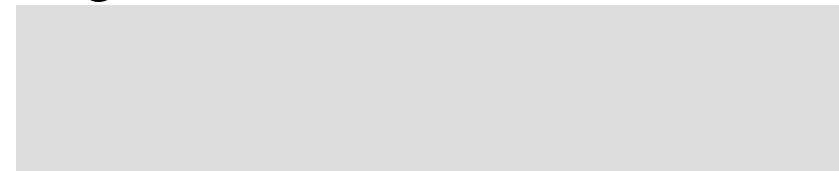
$$\begin{aligned} w_{t+1} &= w_t + \alpha(w_t - 2w_t)\nabla\hat{v} \\ &= w_t + \alpha(1 - 2)w_t \cdot 2 \end{aligned}$$

Off-policy,  $\gamma=1$

From left to right:



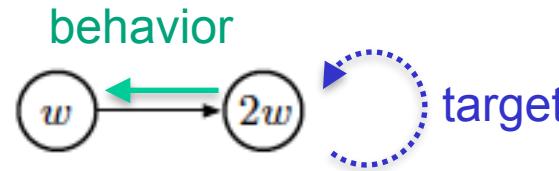
Right to left:



Mean TDE tends to shrink!

# Off-policy prediction with approximation

Initial  $w$  is 10



On-policy,  $\gamma=1$

From left to right:

$$\begin{aligned} w_{t+1} &= w_t + \alpha(2w_t - w_t)\nabla\hat{v} \\ &= w_t + \alpha(2 - 1)w_t \end{aligned}$$

Right to left:

$$\begin{aligned} w_{t+1} &= w_t + \alpha(w_t - 2w_t)\nabla\hat{v} \\ &= w_t + \alpha(1 - 2)w_t \cdot 2 \end{aligned}$$

Mean TDE tends to shrink!

Off-policy,  $\gamma=1$

From left to right:

$$\begin{aligned} w_{t+1} &= w_t + \alpha\rho(2w_t - w_t)\nabla\hat{v} \\ &= w_t + \alpha(2 - 1)w_t \end{aligned}$$

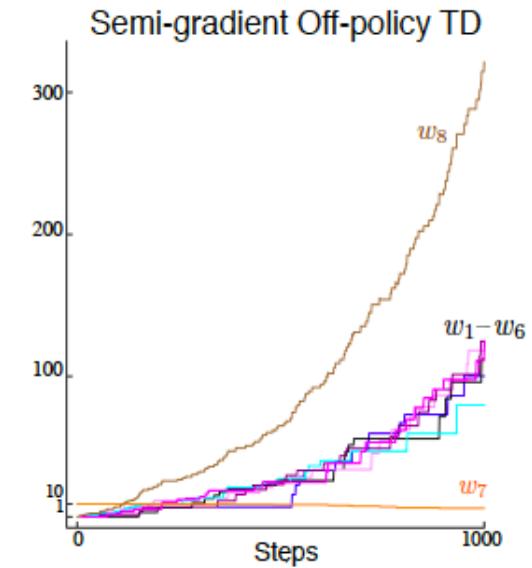
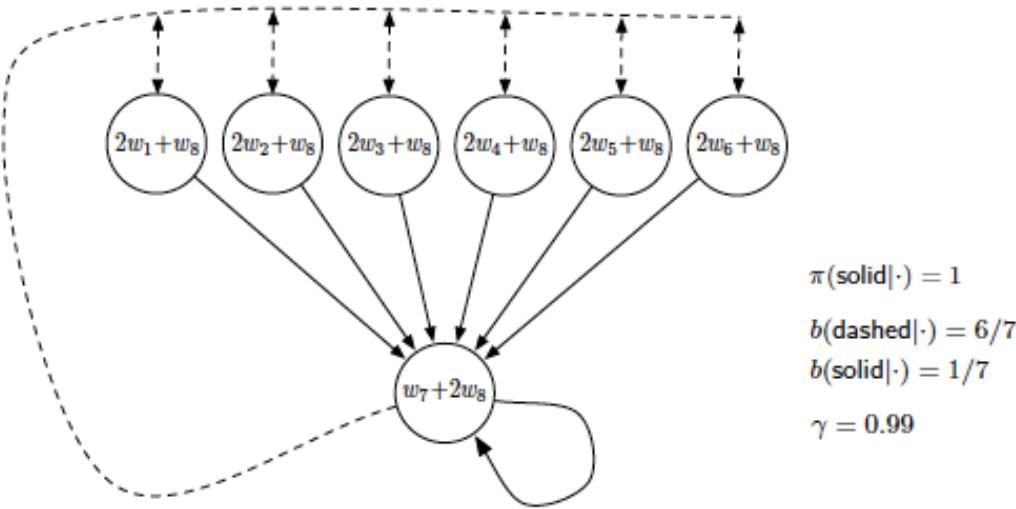
Right to left:

$$w_{t+1} = w_t + \alpha\rho(w - 2w)\nabla\hat{v}$$

If  $\rho \ll 1$ , could TDE keep increasing?

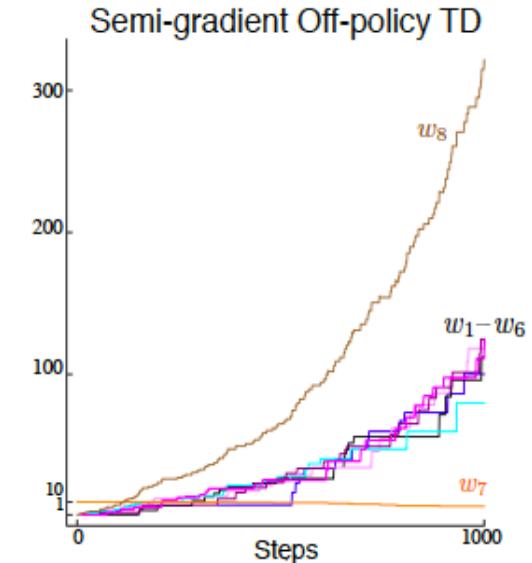
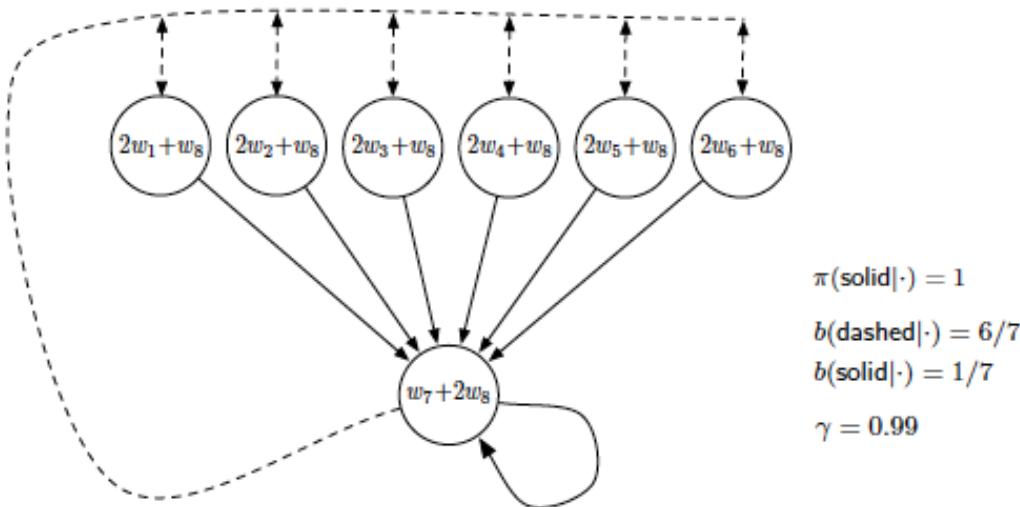
# Off-policy prediction with approximation

Yes, the TDE can keep increasing



# Off-policy prediction with approximation

Yes, the TDE can keep increasing

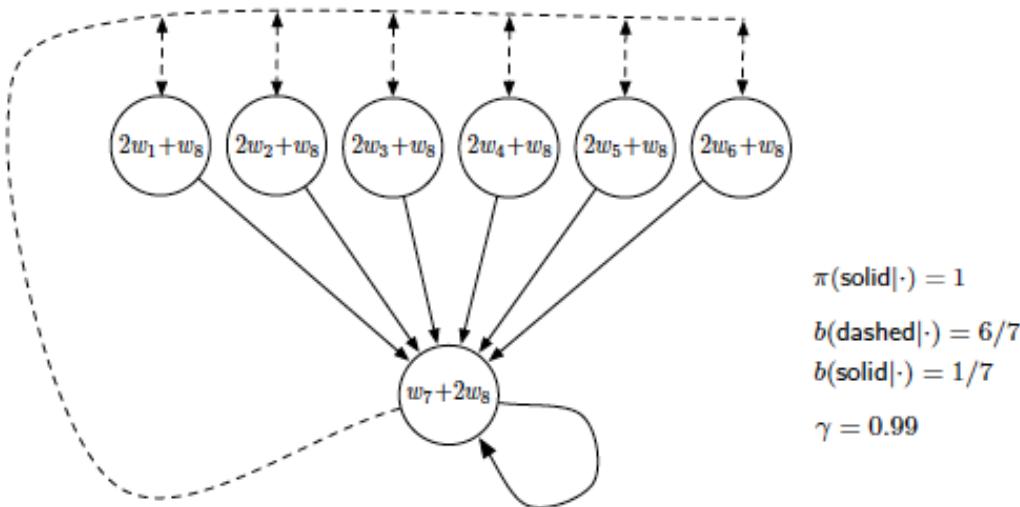


Effect is not because

- $V$  not representable
- Random effects
- Dependent features

# Off-policy prediction with approximation

Yes, the TDE can keep increasing

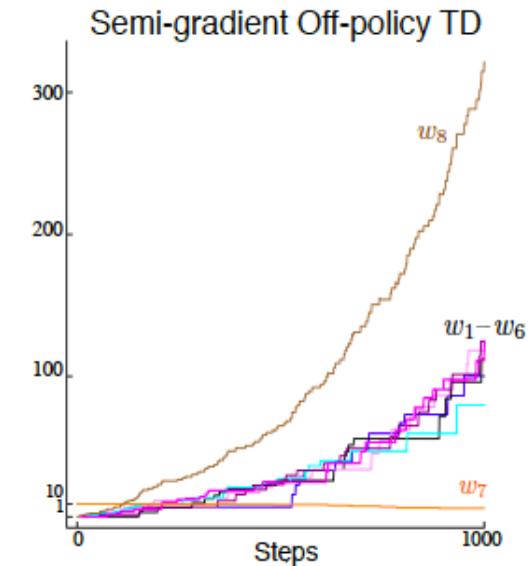


Effect is not because

- V not representable
- Random effects
- Dependent features

But ‘deadly triad’:

- Function approximation
- *Semi-gradient* bootstrapping
- Off-policy training



---

# Off-policy prediction with approximation

---

Can we make semi-gradients work?

Get rid of function approximation?

- Cannot scale to large / continuous domains

Get rid of bootstrapping?

- MC is much slower!

Get rid of off-policy learning?

- Possible, but...
- We would like to learn different policies at same time  
(e.g. running, walking, standing, all from same behavior)
- Also: reuse old data, learn from recorded data, ...

# Alternatives to semi-gradients?

3 states

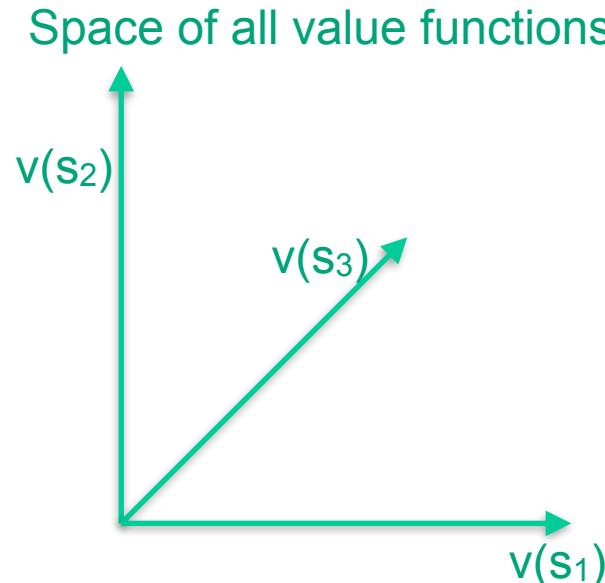
True value

$(v(s_1), v(s_2), v(s_3))$

represented as point  
in 3D space

Approximated value  
 $\hat{v} = w_1x_1(s) + w_2x_2(s)$

- $x(s_1) = [1, 0]$
- $x(s_2) = [-0.5, 0]$
- $x(s_3) = [0, 1]$



Represent as point in 2D subspace

# Alternatives to semi-gradients?

3 states

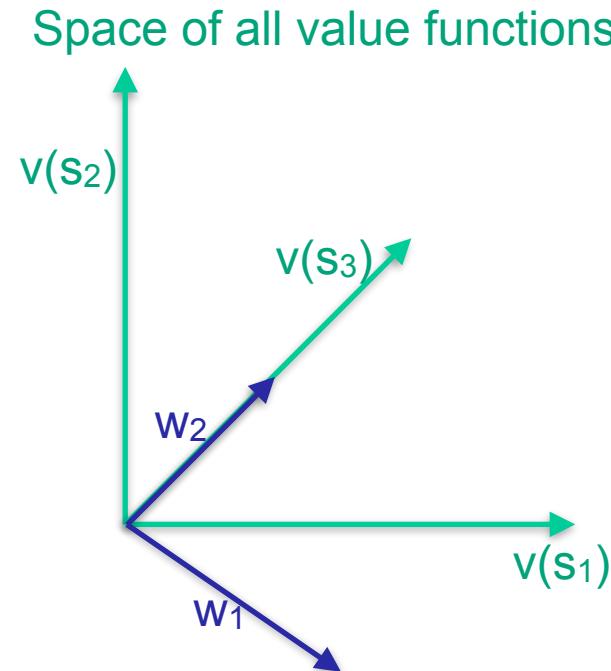
True value

$(v(s_1), v(s_2), v(s_3))$

represented as point  
in 3D space

Approximated value  
 $\hat{v} = w_1x_1(s) + w_2x_2(s)$

- $x(s_1)=[1,0]$
- $x(s_2)=[-0.5,0]$
- $x(s_3)=[0,1]$



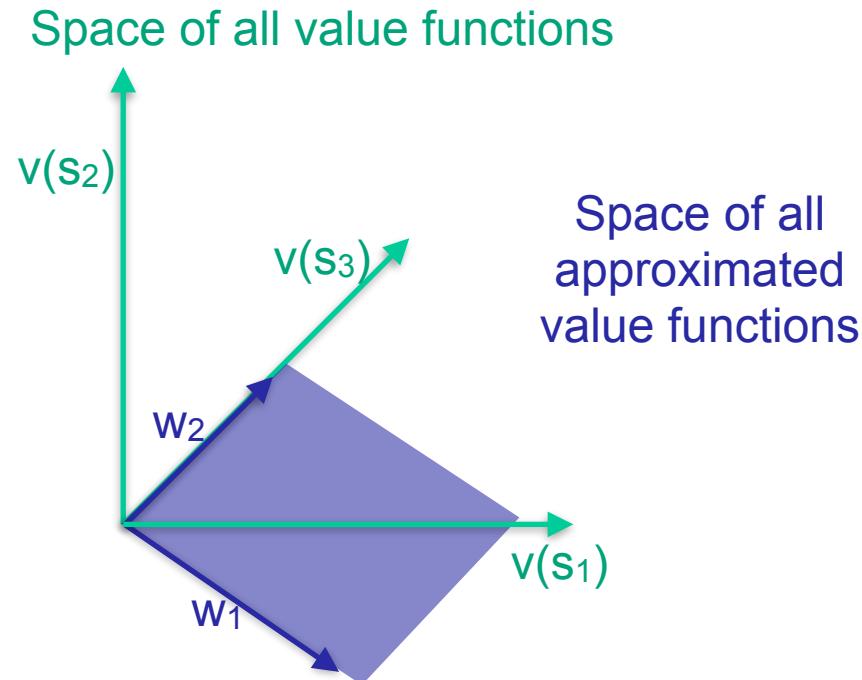
# Alternatives to semi-gradients?

3 states

True value  
 $(v(s_1), v(s_2), v(s_3))$   
represented as point  
in 3D space

Approximated value  
 $\hat{v} = w_1 x_1(s) + w_2 x_2(s)$

- $x(s_1)=[1,0]$
- $x(s_2)=[-0.5,0]$
- $x(s_3)=[0,1]$



Represent as point in 2D subspace

# Alternatives to semi-gradients?

Without approximation, updates approximate Bellman operator

$$B_\pi v_{\mathbf{w}} = v_{\mathbf{w}} + \bar{\delta}_{\mathbf{w}}$$

$$\bar{\delta}_{\mathbf{w}}(s) = \mathbb{E}_\pi \underbrace{[R_{t+1} + \gamma v_{\mathbf{w}}(S_{t+1}) - v_{\mathbf{w}}(S_t)]}_{\delta_{\mathbf{w}}(S_t, A_t, S_{t+1})} | S_t = s, A_t \sim \pi$$

Figure: Sutton & Barto. RL:AI

# Alternatives to semi-gradients?

Without approximation, updates approximate Bellman operator

$$B_\pi v_w = v_w + \bar{\delta}_w$$

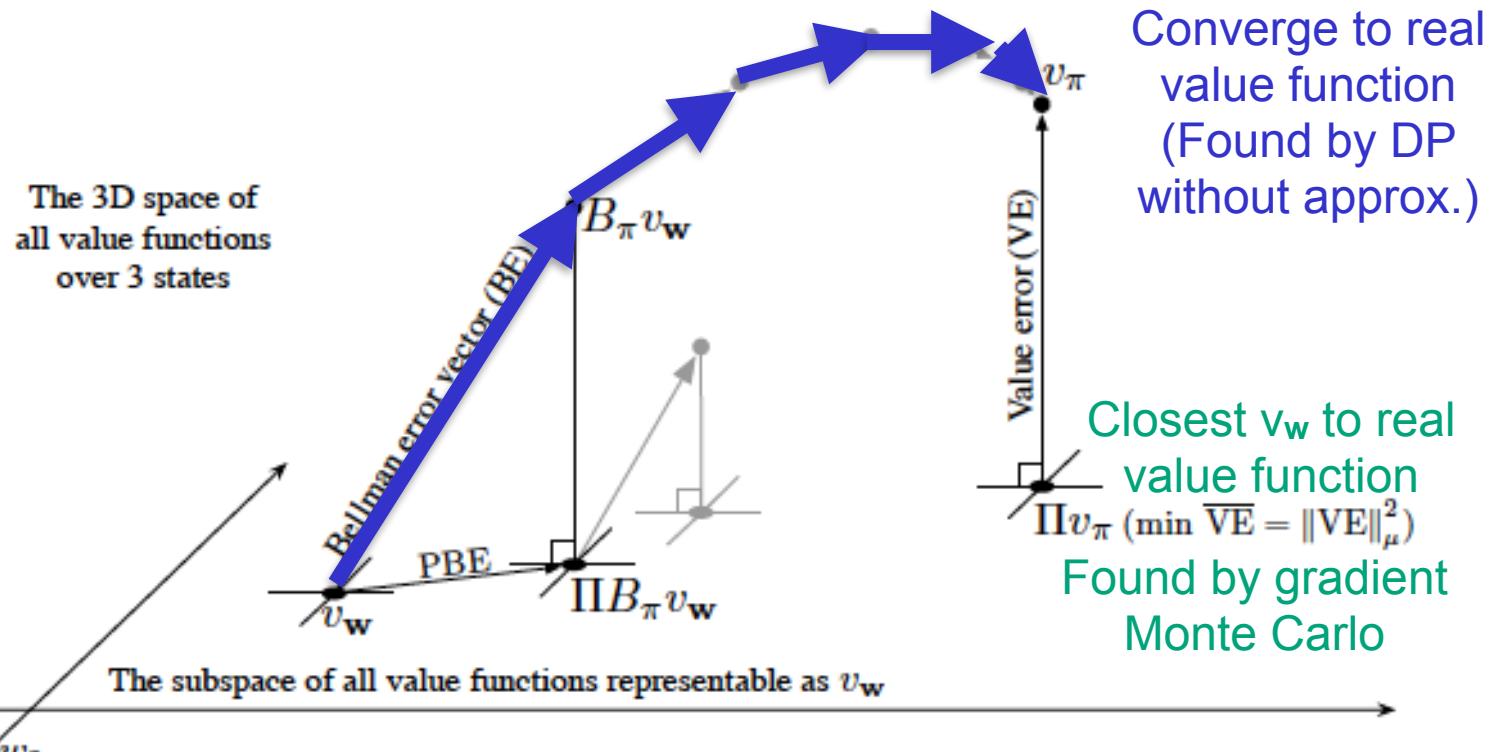


Figure: Sutton & Barto. RL:AI

# Alternatives to semi-gradients?

Without approximation, updates approximate Bellman operator

$$B_\pi v_w = v_w + \bar{\delta}_w$$

The 3D space of all value functions over 3 states

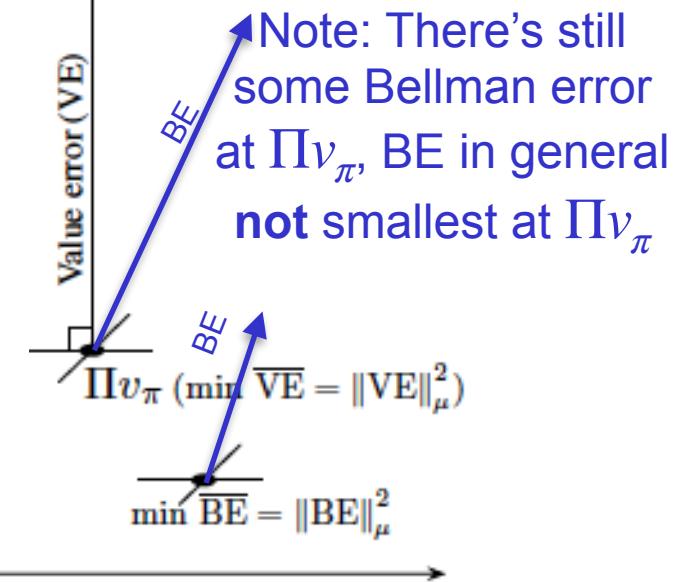
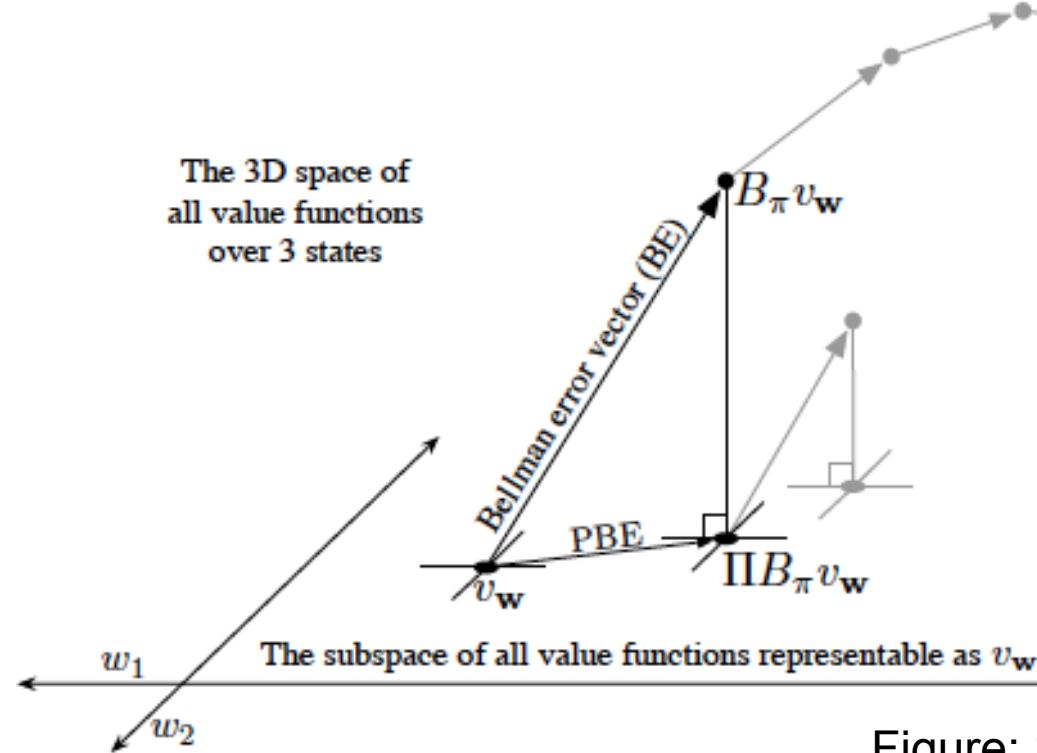


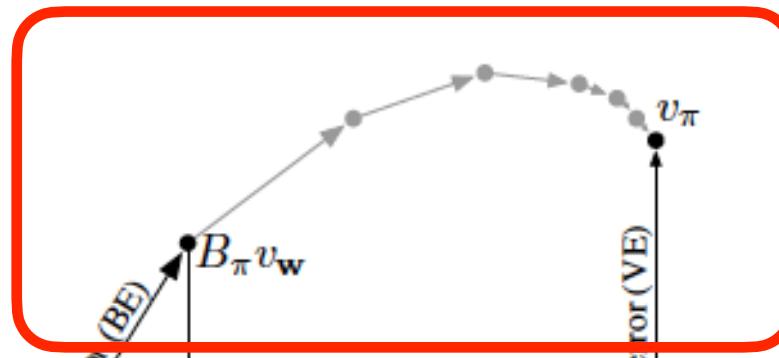
Figure: Sutton & Barto. RL:AI

# Alternatives to semi-gradients?

Without approximation, updates approximate Bellman operator

$$B_\pi v_w = v_w + \bar{\delta}_w$$

The 3D space of all value functions over 3 states



approximation:  
not representable!

$$\Pi v_\pi \quad (\min \overline{VE} = \|VE\|_\mu^2)$$

$$\min \overline{BE} = \|BE\|_\mu^2$$

The subspace of all value functions representable as  $v_w$

$w_1$

$w_2$

Figure: Sutton & Barto. RL:AI

# Alternatives to semi-gradients?

Update according to Bellman error?

$$B_\pi v_w = v_w + \bar{\delta}_w$$

$\Pi$ : Closest alternative  
to desired update

In norm:

$$\|\Delta v\|_\mu^2 \doteq \sum_{s \in S} \mu(s) \Delta v(s)^2$$

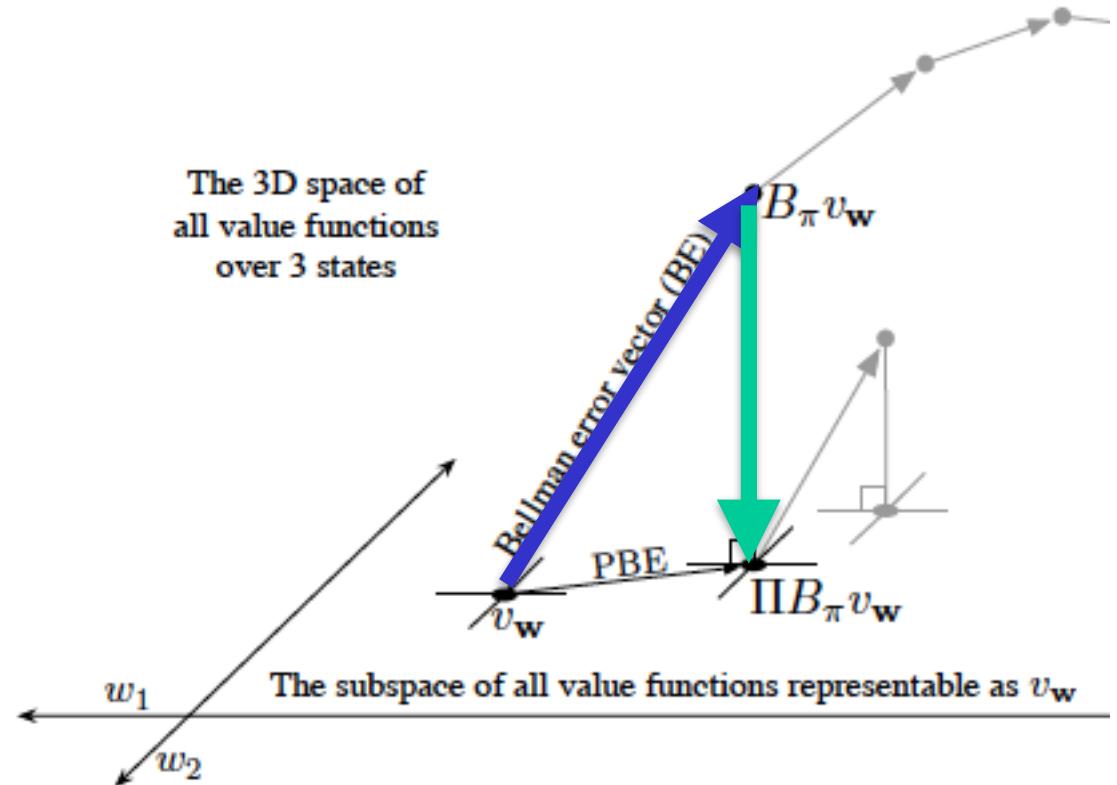
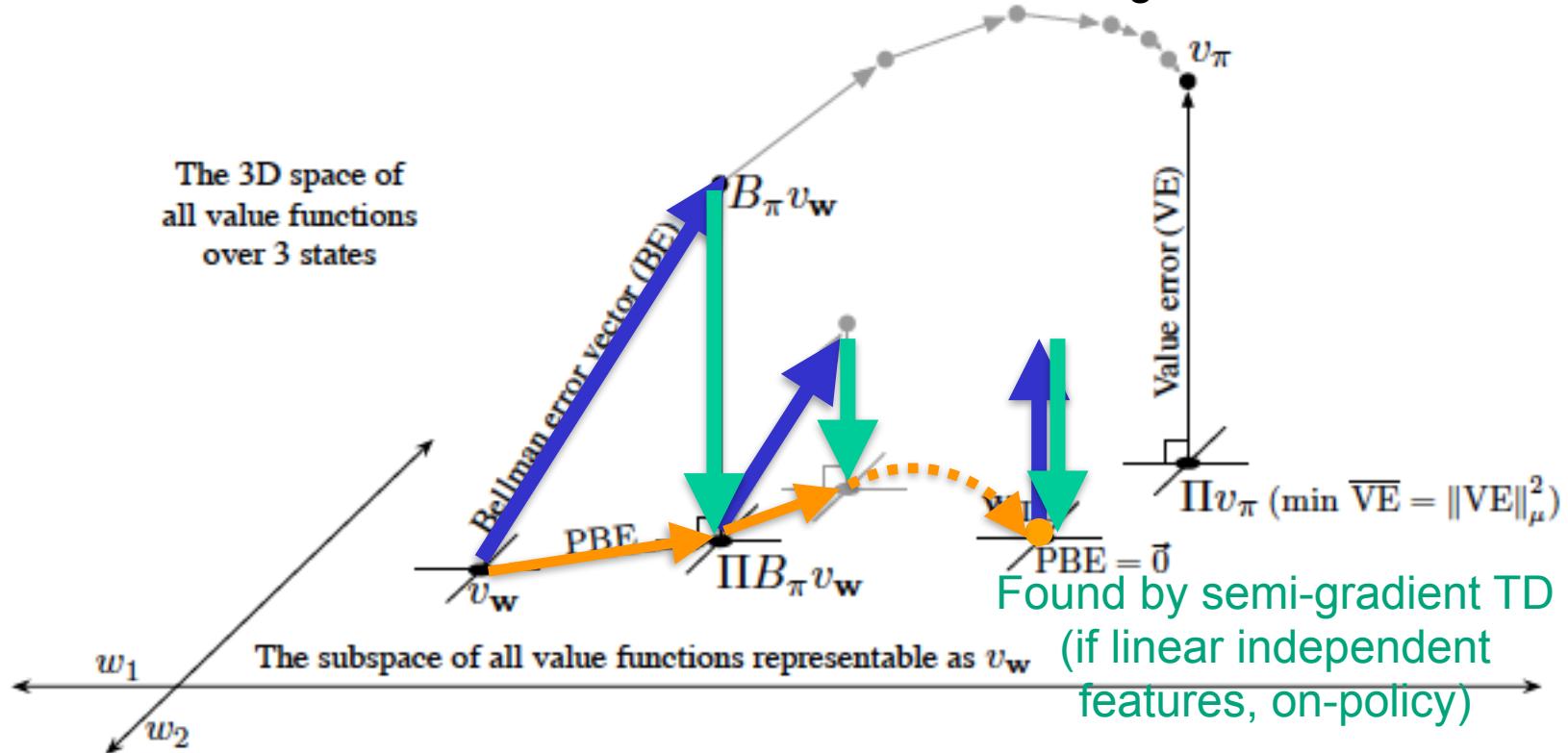


Figure: Sutton & Barto. RL:AI

# Alternatives to semi-gradients?

With approximation, need to stay within representable space

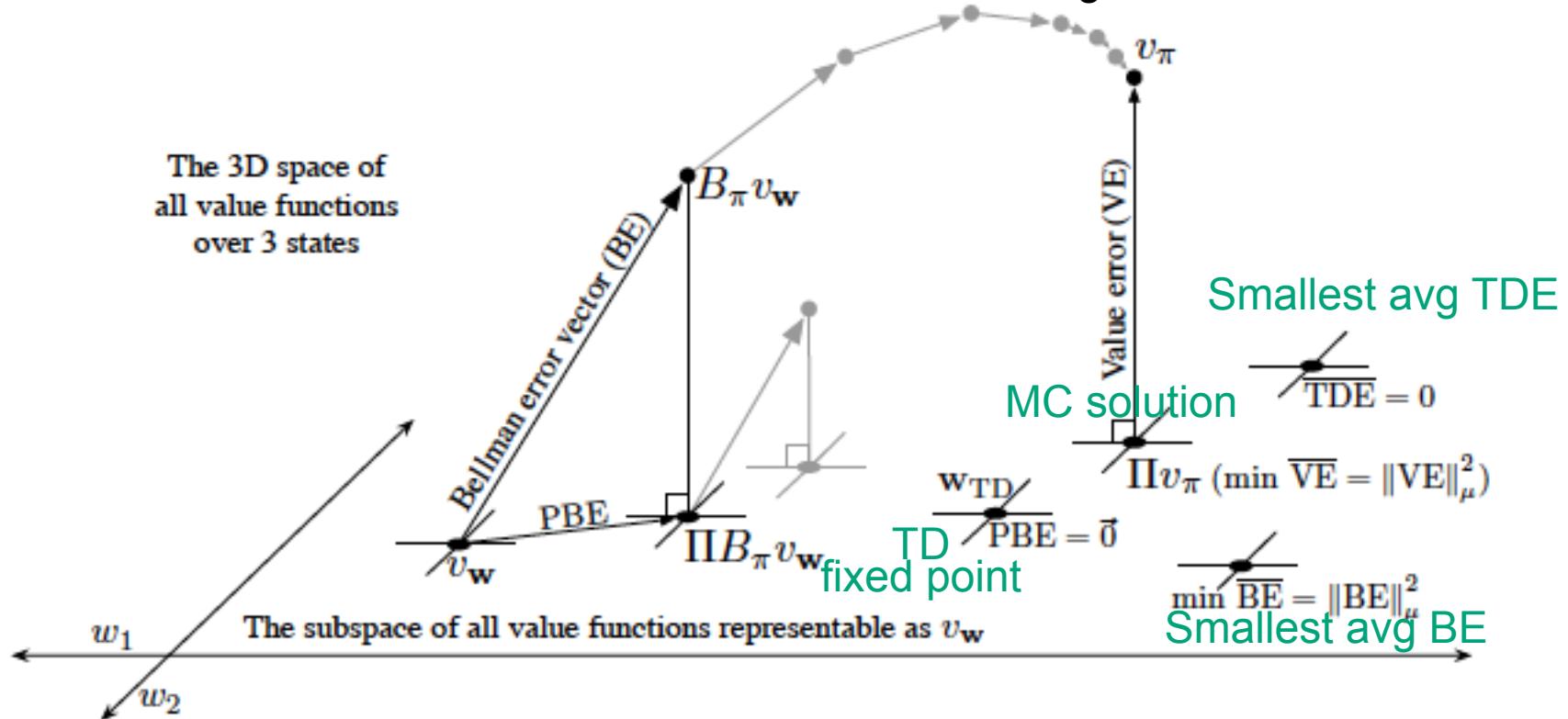
Figure: Sutton & Barto. RL:AI



# Alternatives to semi-gradients?

With approximation, need to stay within representable space

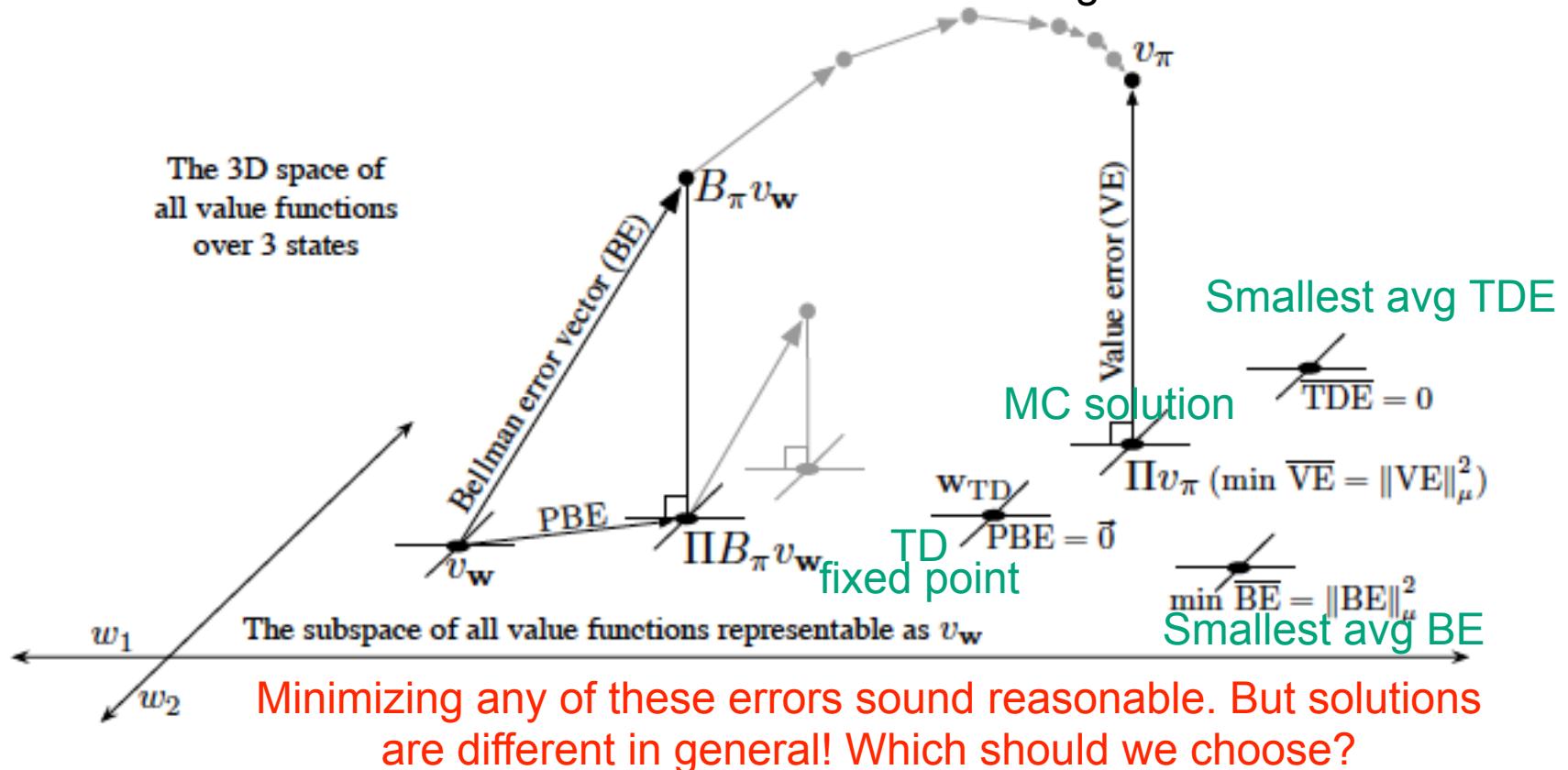
Figure: Sutton & Barto. RL:AI



# Alternatives to semi-gradients?

With approximation, need to stay within representable space

Figure: Sutton & Barto. RL:AI



# Alternatives to semi-gradients?

Find a true gradient algorithm based on minimizing an error?

- Mean squared TD error?

$$\delta_t = R_{t+1} + \gamma v_{\mathbf{w}}(S_{t+1}) - v_{\mathbf{w}}(S_t)$$

$$\overline{\text{TDE}}(\mathbf{w}) = \sum_{s \in S} \mu(s) \mathbb{E} [\delta_t^2 | S_t = s, A_t \sim \pi]$$

Average under visitation distribution  $\mu$

- MS Bellman error?

$$\overline{\text{BE}}(\mathbf{w}) = \sum_{s \in S} \mu(s) \mathbb{E} [\delta_t | S_t = s, A_t \sim \pi]^2$$

- MS Projected Bellman error  $\overline{\text{PBE}}(\mathbf{w}) = \sum_{s \in S} \mu(s) (\Pi \mathbb{E} [\delta_t | S_t = s, A_t \sim \pi])^2$

# Minimize mean squared TD error?

---

$$\begin{aligned}\overline{\text{TDE}}(\mathbf{w}) &= \sum_{s \in S} \mu(s) \mathbb{E} [\delta_t^2 | S_t = s, A_t \sim \pi] \\ &= \sum_{s \in S} \mu(s) \mathbb{E} [\rho_t \delta_t^2 | S_t = s, A_t \sim b] \\ &= \mathbb{E}_b [\rho_t \delta_t^2].\end{aligned}$$

We can estimate the real gradient of the TDE from data!  
How does it perform?

# Minimize mean squared TD error?

## A-split example

- What value minimizes TD error at A?
- Now, what value minimizes TD error at B, C?

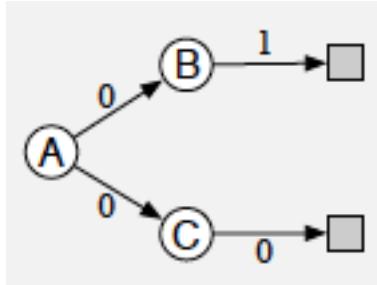


Figure: Sutton & Barto. RL:AI

# Minimize mean squared Bellman error?

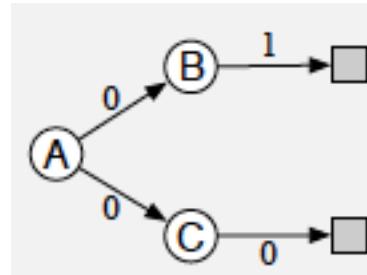
Let's try something else

$$\overline{\text{BE}}(\mathbf{w}) = \|\bar{\delta}_{\mathbf{w}}\|_{\mu}^2$$

$$\bar{\delta}_{\mathbf{w}}(s) = \mathbb{E}_{\pi} \underbrace{[R_{t+1} + \gamma v_{\mathbf{w}}(S_{t+1}) - v_{\mathbf{w}}(S_t)]}_{\delta_{\mathbf{w}}(S_t, A_t, S_{t+1})} | S_t = s, A_t \sim \pi$$

Least Bellman error is acceptable solution in A-split problem

However, the gradient cannot be estimated from data only!



# Mean squared projected Bellman error?

PBE last chance for bootstrapping!

- Linear approximation: PBE=0 at  $w_{TD}$
- On-policy, semi-gradient finds  $w_{TD}$
- LSTD find TD fixed point
- Find TD fixed point off-policy?

Let's minimize PBE used gradient descent algorithm

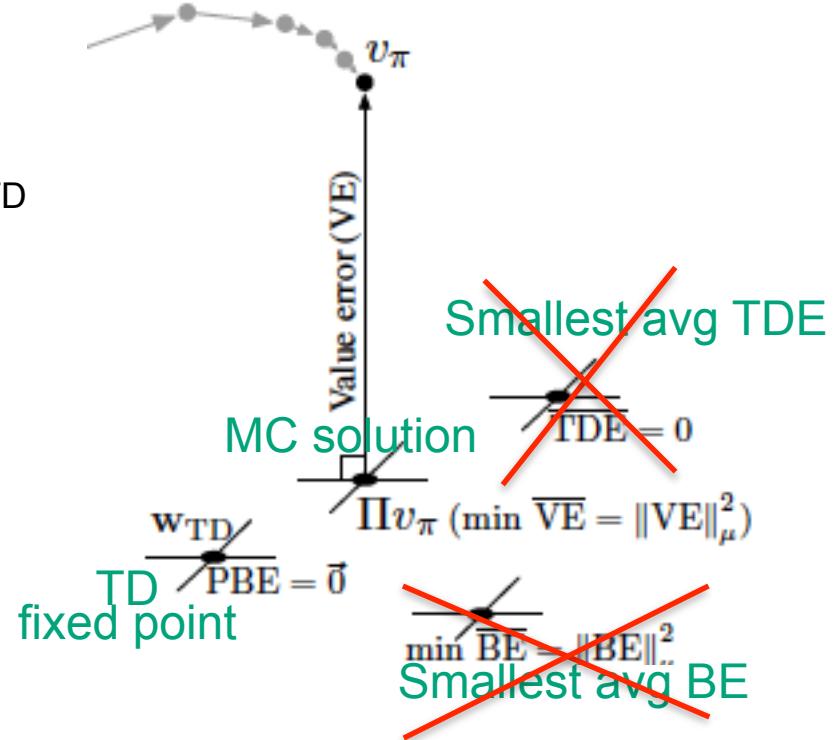


Figure: Sutton & Barto. RL:AI

# Mean squared projected Bellman error?

---

It can be shown that

$$\begin{aligned}\nabla \overline{\text{PBE}}(\mathbf{w}) &= \nabla \left\| \Pi \bar{\delta}_{\mathbf{w}} \right\|_{\mu}^2 \\ &= 2\mathbb{E} \left[ \rho_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t) \mathbf{x}_t^\top \right] \mathbb{E} \left[ \mathbf{x}_t \mathbf{x}_t^\top \right]^{-1} \mathbb{E} [\rho_t \delta_t \mathbf{x}_t]\end{aligned}$$

Reproducing the exact steps is outside the scope of the course  
(the book gives a good explanation if you're curious)

# Mean squared projected Bellman error?

---

It can be shown that

$$\begin{aligned}\nabla \overline{\text{PBE}}(\mathbf{w}) &= \nabla \left\| \Pi \bar{\delta}_{\mathbf{w}} \right\|_{\mu}^2 \\ &= 2\mathbb{E} \left[ \rho_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t) \mathbf{x}_t^\top \right] \mathbb{E} \left[ \mathbf{x}_t \mathbf{x}_t^\top \right]^{-1} \mathbb{E} [\rho_t \delta_t \mathbf{x}_t]\end{aligned}$$

Reproducing the exact steps is outside the scope of the course  
(the book gives a good explanation if you're curious)

**Can we just plug in a new samples ( $s,a,s'$ ) tuple in each  $\mathbb{E}$ ?**

# Mean squared projected Bellman error?

$$\nabla \overline{\text{PBE}}(\mathbf{w}) = 2\mathbb{E} [\rho_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t) \mathbf{x}_t^\top] \mathbb{E} [\mathbf{x}_t \mathbf{x}_t^\top]^{-1} \mathbb{E} [\rho_t \delta_t \mathbf{x}_t]$$

The first and last factor are dependent - both depend on  $\mathbf{x}_{t+1}$

Why a problem? Example: consider r.v.  $Y \in \{-1, 1\}$  uniform

$$\mathbb{E}[Y]\mathbb{E}[Y]$$

$$=?$$

$$\sum_{i=1}^N Y_i Y_i$$

$$=?$$

# Mean squared projected Bellman error?

---

$$\nabla \overline{\text{PBE}}(\mathbf{w}) = 2\mathbb{E} [\rho_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t) \mathbf{x}_t^\top] \mathbb{E} [\mathbf{x}_t \mathbf{x}_t^\top]^{-1} \mathbb{E} [\rho_t \delta_t \mathbf{x}_t]$$

The first and last factor are dependent - both depend on  $\mathbf{x}_{t+1}$   
Learn some factors from all data, only plug in sample for rest

# Mean squared projected Bellman error?

$$\nabla \overline{\text{PBE}}(\mathbf{w}) = 2\mathbb{E} [\rho_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t) \mathbf{x}_t^\top] \mathbb{E} [\mathbf{x}_t \mathbf{x}_t^\top]^{-1} \mathbb{E} [\rho_t \delta_t \mathbf{x}_t]$$

---

v

The first and last factor are dependent - both depend on  $\mathbf{x}_{t+1}$   
Learn some factors from all data, only plug in sample for rest

# Mean squared projected Bellman error?

$$\nabla \overline{\text{PBE}}(\mathbf{w}) = 2\mathbb{E} [\rho_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t) \mathbf{x}_t^\top] \mathbb{E} [\mathbf{x}_t \mathbf{x}_t^\top]^{-1} \mathbb{E} [\rho_t \delta_t \mathbf{x}_t]$$

---

v

The first and last factor are dependent - both depend on  $\mathbf{x}_{t+1}$   
Learn some factors from all data, only plug in sample for rest  
 $\mathbf{v}$  is the solution to a least squares problem - for such problems  
there is a SGD algorithm:  $\mathbf{v}_{t+1} \doteq \mathbf{v}_t + \beta \rho_t (\delta_t - \mathbf{v}_t^\top \mathbf{x}_t) \mathbf{x}_t$

Plug in the parts:  $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha / 2 \nabla \overline{\text{PBE}}(\mathbf{w}_t)$

$$\begin{aligned} &\approx \mathbf{w}_t + \alpha \mathbb{E} [\rho_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1}) \mathbf{x}_t^\top] \mathbf{v}_t \\ &\approx \mathbf{w}_t + \alpha \rho_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1}) \mathbf{x}_t^\top \mathbf{v}_t \end{aligned}$$

# Mean squared projected Bellman error?

---

This is GTD2, a gradient-TD method

- Converges to 0 mean squared projected Bellman error (linear features)
- With non-linear feature approximation, converges to local optimum
- There are other, more sophisticated gradient-TD methods
  
- Cost?
- Additional step size parameter for learning  $\mathbf{v}$  ('larger' than  $\alpha$ )
- Need to store and update two parameter vectors
- Extension needed to apply it to non-linear approximators
- *If* semi-gradient TD converges, it seems to converge faster in practice

# Convergence with function approximation

	Tabular On/Off	Linear on **	Nonlinear on	Linear off **	Nonlinear off
Gradient MC *					
Semi-gradient TD *				No C!	No C!
Gradient TD *					
LSTD				N.A.	N.A.

\* with appropriate step-size schedule

\*\* if features independent, single solution

# Convergence with function approximation

	Tabular On/Off	Linear on **	Nonlinear on	Linear off **	Nonlinear off
Gradient MC *					
Semi-gradient TD *			Convergence to minimum of which error? !?	No C!	No C!
Gradient TD *					
LSTD				N.A.	N.A.

\* with appropriate step-size schedule

\*\* if features independent, single solution

# Convergence with function approximation

	Tabular On/Off	Linear on **	Nonlinear on	Linear off **	Nonlinear off
Gradient MC *	MC:VE				
Semi-gradient TD *	TD: PBE				
			No C!	No C!	No C!
Gradient TD *	TD: PBE				
LSTD	TD: PBE				
			N.A.		N.A.

\* with appropriate step-size schedule

\*\* if features independent, single solution

# Convergence with function approximation

	Tabular On/Off	Linear on **	Nonlinear on	Linear off **	Nonlinear off
Gradient MC *					
Semi-gradient TD *			Local or global convergence?	No C!	No C!
Gradient TD *					
LSTD				N.A.	N.A.

\* with appropriate step-size schedule

\*\* if features independent, single solution

# Convergence with function approximation

	Tabular On/Off	Linear on **	Nonlinear on	Linear off **	Nonlinear off	
Gradient MC *	global	global	non-linear: local	No C!	global	non-linear: local
Semi-gradient TD *	global	global	non-linear: local	No C!	global	No C!
Gradient TD *	global	global	non-linear: local	No C!	global	non-linear: local
LSTD			N.A.			N.A.

\* with appropriate step-size schedule

\*\* if features independent, single solution

# Convergence with function approximation

	Tabular On/Off	Linear on **	Nonlinear on	Linear off **	Nonlinear off
Gradient MC *	global	global	MC:VE		
Semi-gradient TD *	global	global	TD: PBE	No C!	No C!
Gradient TD *	global	non-linear: local	TD: PBE	global	non-linear: local
LSTD			TD: PBE	N.A.	N.A.

\* with appropriate step-size schedule

\*\* if features independent, single solution

# Deep Q networks

The famous algorithm for learning to play  
Atari games: Deep Q networks (DQN)

Illustrates some of the challenges in doing  
control with (non-linear) function  
approximation

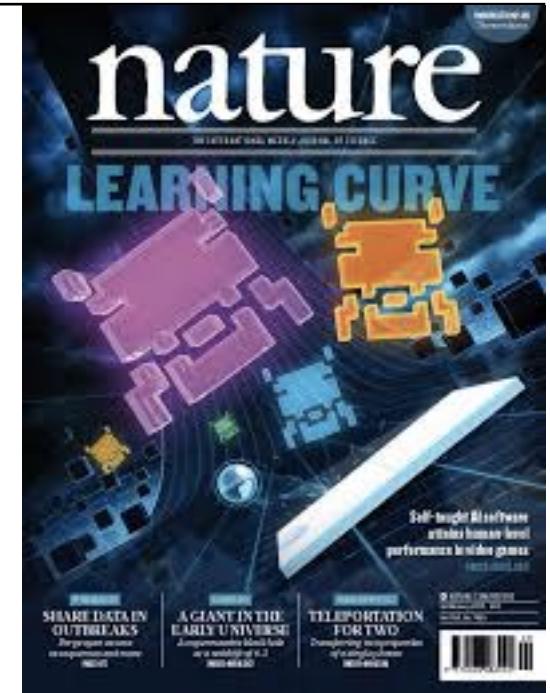
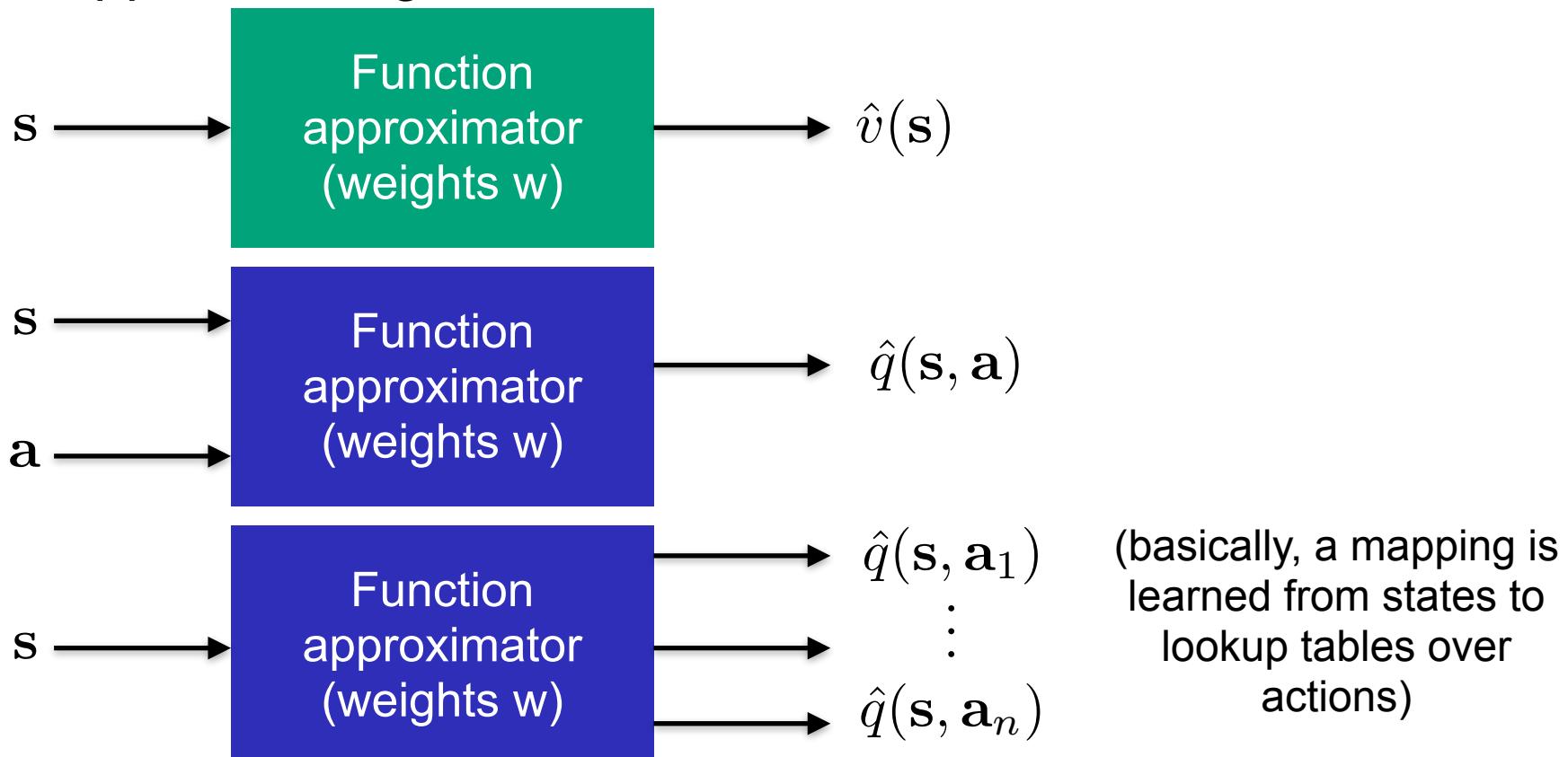


Figure: Deepmind.com

# Deep Q networks

DQN approximates the Q function. Two possibilities compared to approximating V:



# Deep Q networks

---

Linear approximation has nice properties, but task-specific features need to be designed

DQN: Instead, use a multi-layer neural network. Idea is that early layers can learn features

Thus, the network can directly take relatively ‘raw’ input without task-specific preprocessing

# Deep Q networks

## Pre-processing (Task-agnostic )

- Downscale image and convert to grayscale
- Stack frames
- Reward = increment in score



Going left or right?

Figure: Smithsonianmag.com

# Deep Q networks

Most layers in DQN are convolutional layers:

- Same filters are used at each location in the image
- Reduces #weights by only connecting neurons locally
- Further reduces #weights by sharing parameters at each location
- Imposes translation equivariance

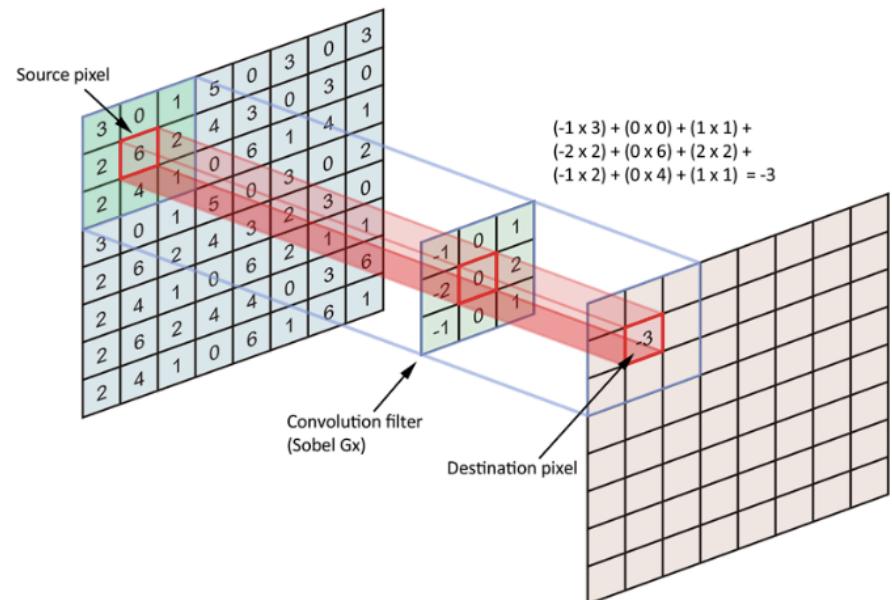


Figure: towardsdatascience.com

# Deep Q networks

Convolutional layers and fully connected layers of ReLU units

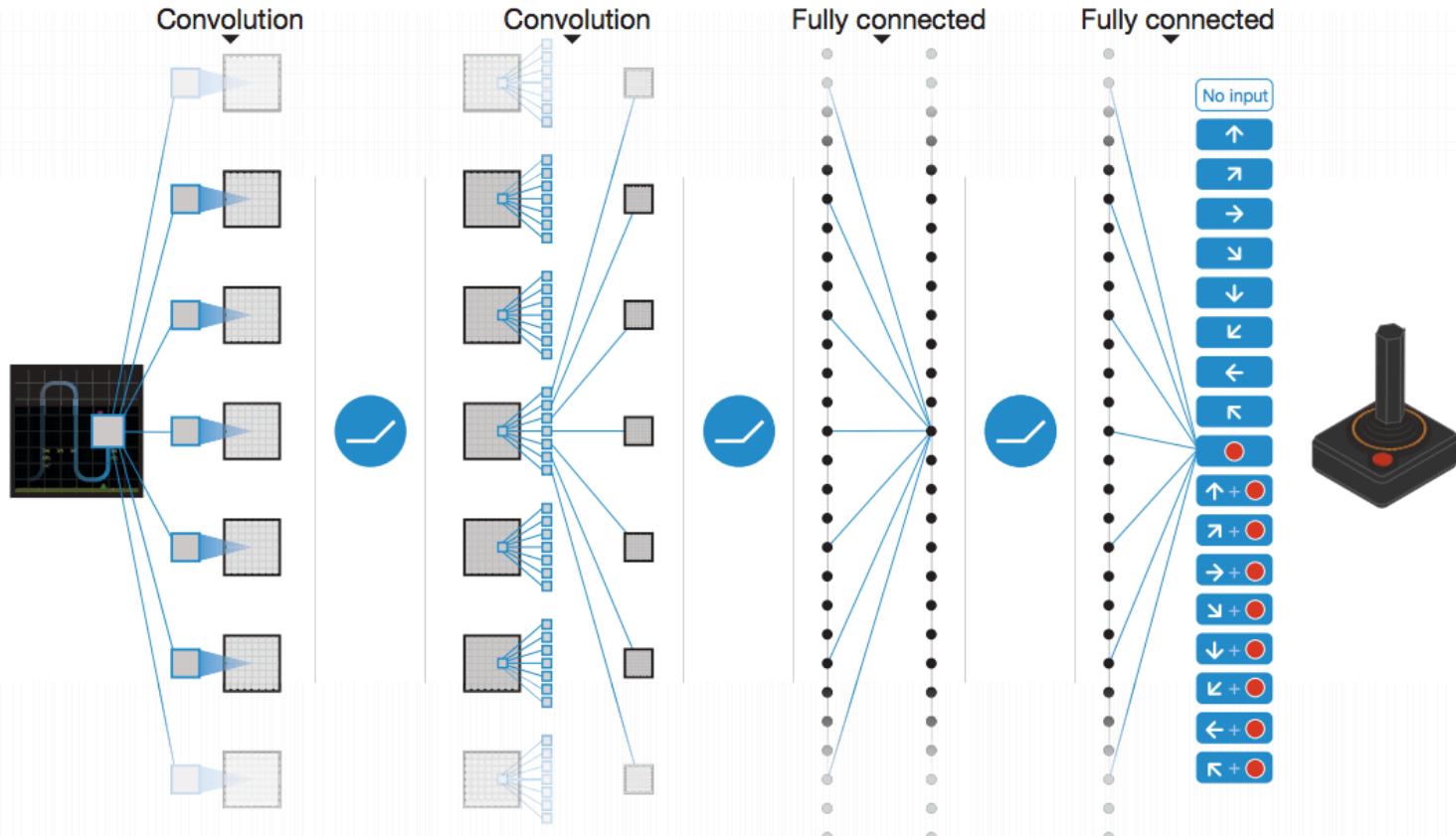


Figure: Mnih et al., 2015

# Deep Q networks

---

In regular TD-learning, each transition used once. Inefficient  
Also, subsequent transitions are highly **correlated**

Break correlation using experience replay

- Store experiences in a buffer
- At every time step, train on a batch of experiences from the buffer, not the most recent one
- Makes sure the network can't 'overfit' to recent experience only

Are we learning on-policy or off-policy?

# Deep Q networks

Learning with semi-gradient version of q-learning

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left[ \underbrace{R + \gamma \max_a \hat{q}(S', a, \mathbf{w}) - \hat{q}(S_t, a, \mathbf{w})}_{\text{Target depends on w}} \right] \nabla \hat{q}(S_t, a, \mathbf{w})$$

Dependence of target on  $\mathbf{w}$  can cause instabilities

- Copy parameters to  $\tilde{\mathbf{w}}$  every C steps, and use these in the target

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left[ R + \gamma \max_a \hat{q}(S', a, \tilde{\mathbf{w}}) - \hat{q}(S_t, a, \mathbf{w}) \right] \nabla \hat{q}(S_t, a, \mathbf{w})$$


# Deep Q networks

Learning with semi-gradient version of q-learning

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left[ \underbrace{R + \gamma \max_a \hat{q}(S', a, \mathbf{w}) - \hat{q}(S_t, a, \mathbf{w})}_{\text{Target depends on w}} \right] \nabla \hat{q}(S_t, a, \mathbf{w})$$

Dependence of target on  $\mathbf{w}$  can cause instabilities

- Copy parameters to  $\tilde{\mathbf{w}}$  every C steps, and use these in the target

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left[ R + \gamma \max_a \hat{q}(S', a, \tilde{\mathbf{w}}) - \hat{q}(S_t, a, \mathbf{w}) \right] \nabla \hat{q}(S_t, a, \mathbf{w})$$


Extra trick: clip error to  $[-1, 1]$  (more stability, easy to set  $\alpha$ )

# Deep Q networks

Results with the **same** parameters and settings on many games

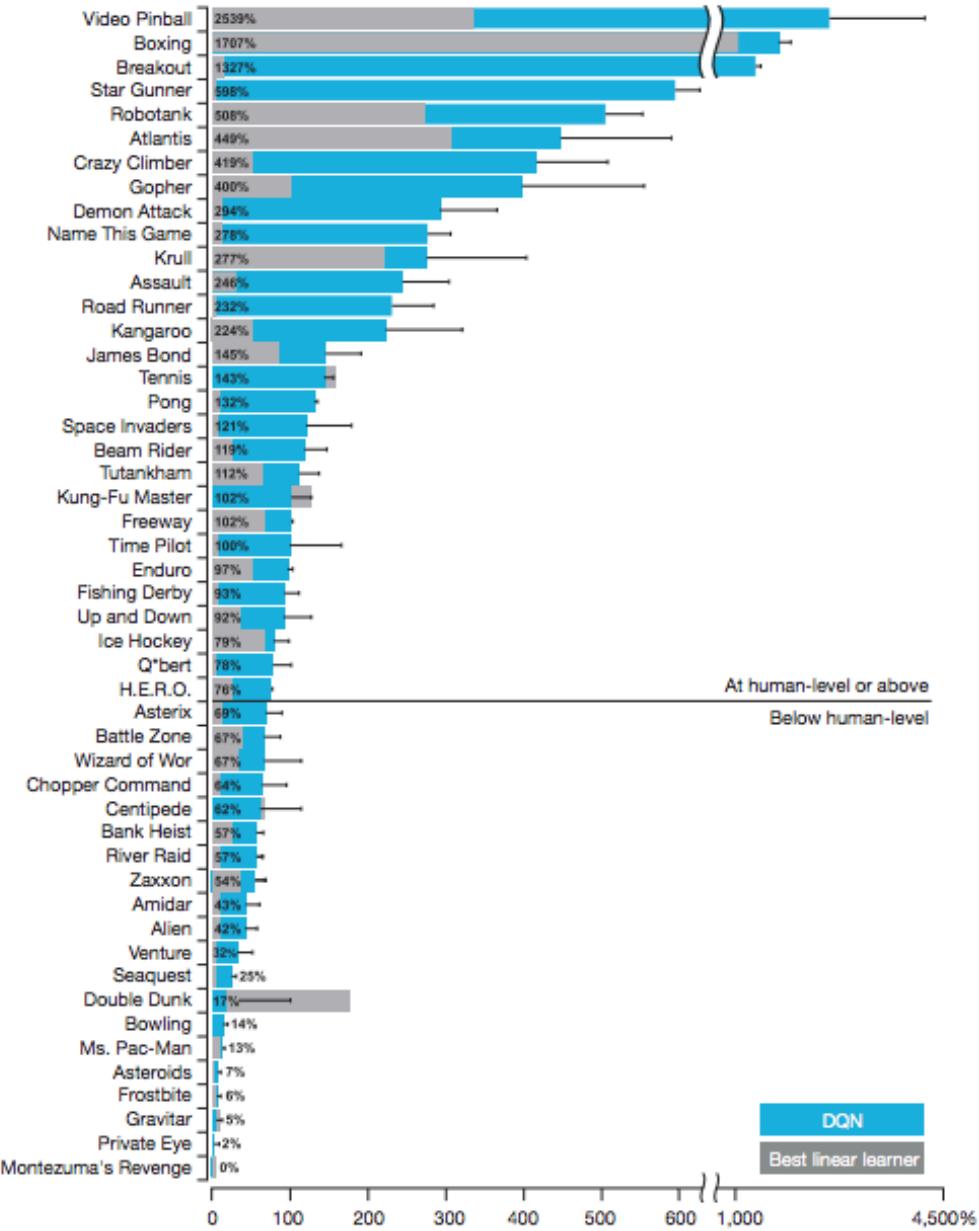


Figure: Mnih et al., 2015

# Deep Q networks

---

Many of the tricks make the problem more ‘like supervised’

- Target network: fixed target value
- Experience replay: closer to iid

Still semi-gradient used off-policy

- Hence the need to stabilise learning
- Lack theoretical guarantees, but good empirical performance

Non-linear function approximation

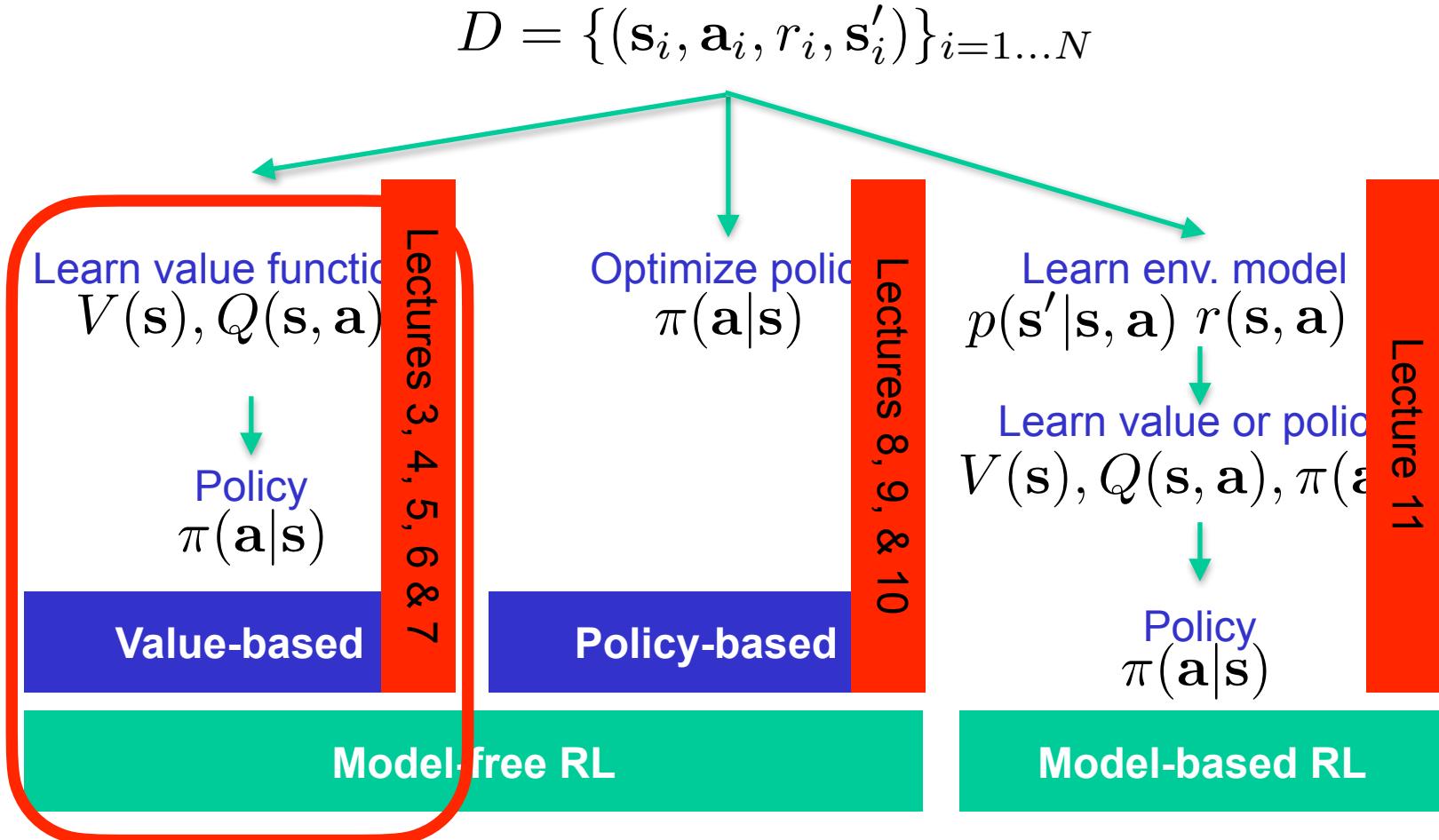
- Risk of local optima seems not so bad in practice for deep nets

Extension: double DQN (similar to double Q learning)

Powerful demonstration of RL without task-specific features!

---

# Big picture: How to learn policies



Thanks to Jan Peters

# Back to the big picture

Last week: *prediction* methods

Approximate method

semi-gradient TD(0)  
& LSTD

Exact methods

gradient  
Monte-Carlo

Temporal-  
difference  
learning

N-Step

Monte  
Carlo

width  
of update

need dynamics  
model



Figure from Sutton and Barto RL:AI  
ng

# Control methods

On-policy  
(narrow)



Sarsa

N-step  
SARSA

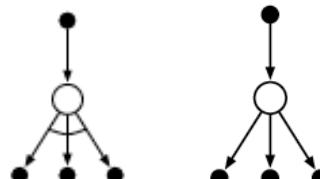


Monte-  
Carlo  
control

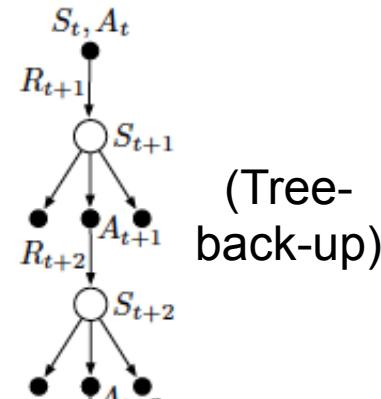


Depth

Off-policy  
(wide)



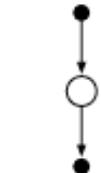
Q-learning &  
Expected SARSA



(Tree-  
back-up)

# Control methods

On-policy  
(narrow)



Sarsa

Semi-gradient  
SARSA



N-step  
SARSA

Depth

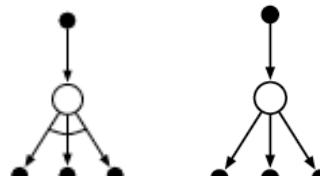
(Semi-  
gradient n-  
step SARSA)



Monte-  
Carlo  
control

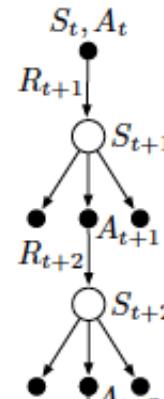
(Gradient  
MC)

Off-policy  
(wide)



Q-learning &  
Expected SARSA

DQN  
GTD2 (off-policy  
but 'narrow')



(Tree-  
back-up)

---

# Conclusion

---

On-policy control straightforward

Off-policy prediction and control very tricky!

- Gradient TD, that uses gradient of PBE, is a solution
- Overhead: keep track of more values
- GTD2 as discussed, only works for linear f.a.:  
not really applicable to ‘raw data’

Alternative: use additional mechanisms to keep learning stable

- Good empirical performance with DQN directly from ‘raw’ data!

---

# Thanks for your attention?

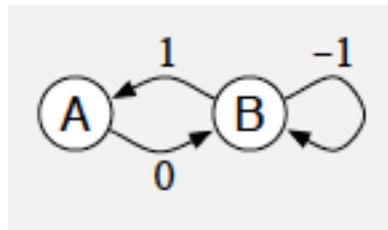
---

Feedback?

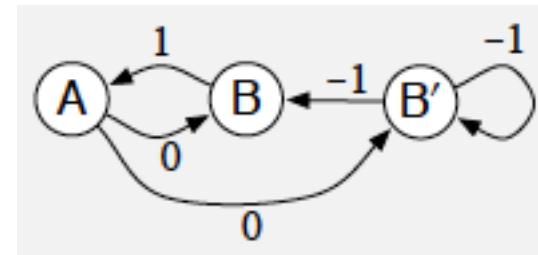
[h.c.vanhoof@uva.nl](mailto:h.c.vanhoof@uva.nl)

# Minimize mean squared Bellman error?

Is it really impossible to learn from data?



$$\begin{aligned}x(A) &= [1, 0]^T \\x(B) &= [0, 1]^T \\x(B') &= [0, 1]^T\end{aligned}$$



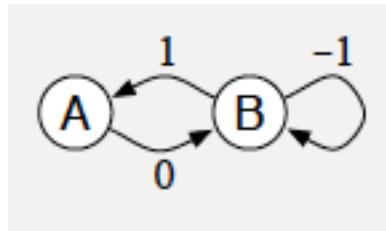
Observed data?

Best  $w$ ?

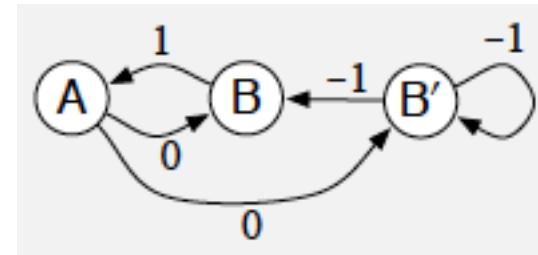
Figure: Sutton & Barto. RL:AI

# Minimize mean squared Bellman error?

Is it really impossible to learn from data?



$$\begin{aligned}x(A) &= [1, 0]^T \\x(B) &= [0, 1]^T \\x(B') &= [0, 1]^T\end{aligned}$$



Observed data?

- Same for both models

Best  $w$ ?

- $[0,0]^T$  for model 1, but this doesn't have 0 error for model 2
- In fact, BE minimised at  $[-0.5, 0]^T$  for model 2 ( $\gamma=1$ )

Figure: Sutton & Barto. RL:AI