
Prediction with approximation

Herke van Hoof

Why approximation

So far, we have looked at ‘small’ problems: can store the Q/V fc for every state

With a huge number of states, two problems:

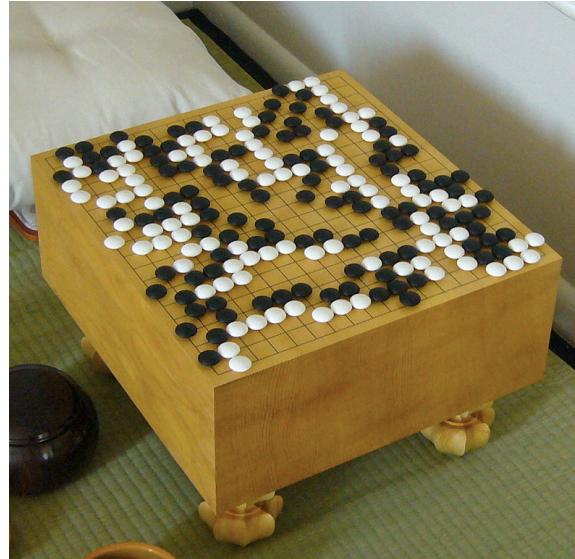
- The table might be too large to fit in memory
- Before memory becomes a problem, it might take too long to collect enough experiences to fill up every row in the table
- The state space might be continuous (infinitely many states)

Today, we’ll look at how to overcome these problems in the *policy evaluation* (i.e., prediction) setting. *Control* next week!

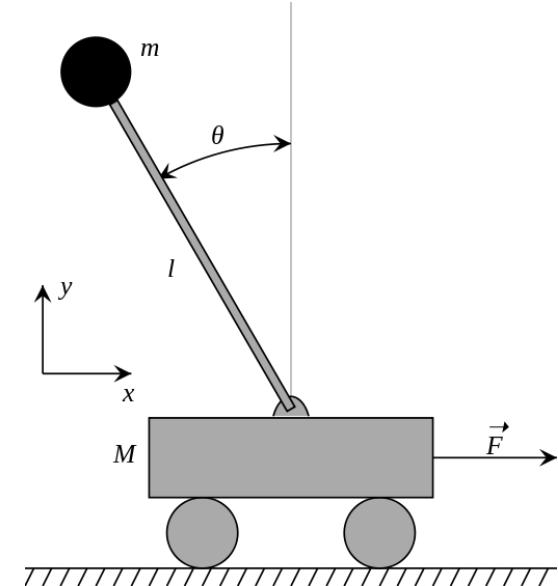
Why approximation



Atari games



Go

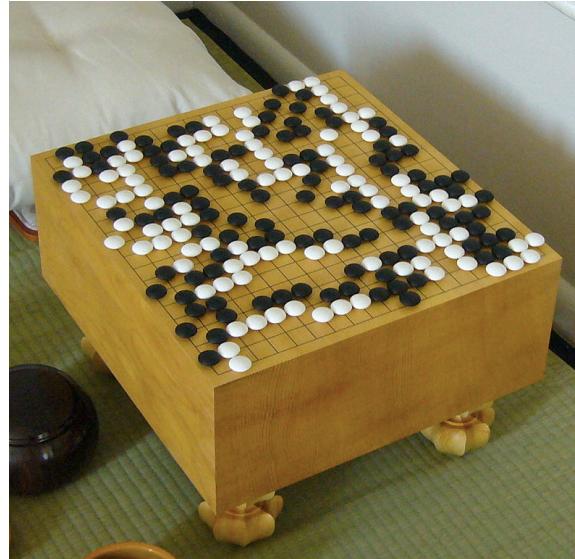


Physical systems

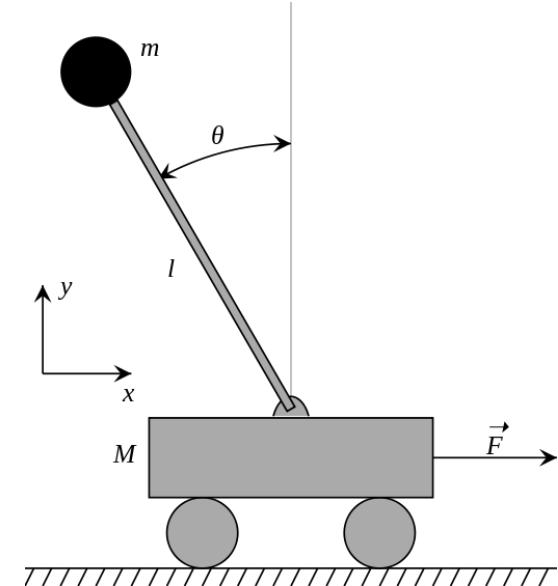
Why approximation



Atari games



Go



Physical systems

Luckily, certain small changes in state usually don't impact the value/action too much. **We can generalize.**

Why approximation

So we would like to represent the value function

- In a compact way
- In a way that allows generalizing an experience to nearby states

This is similar to the goal of supervised learning

- Very flexible function might not generalise well (overfitting)
- Inflexible functions are not very expressive
- Need to find a good balance, need for expressivity <> data availability

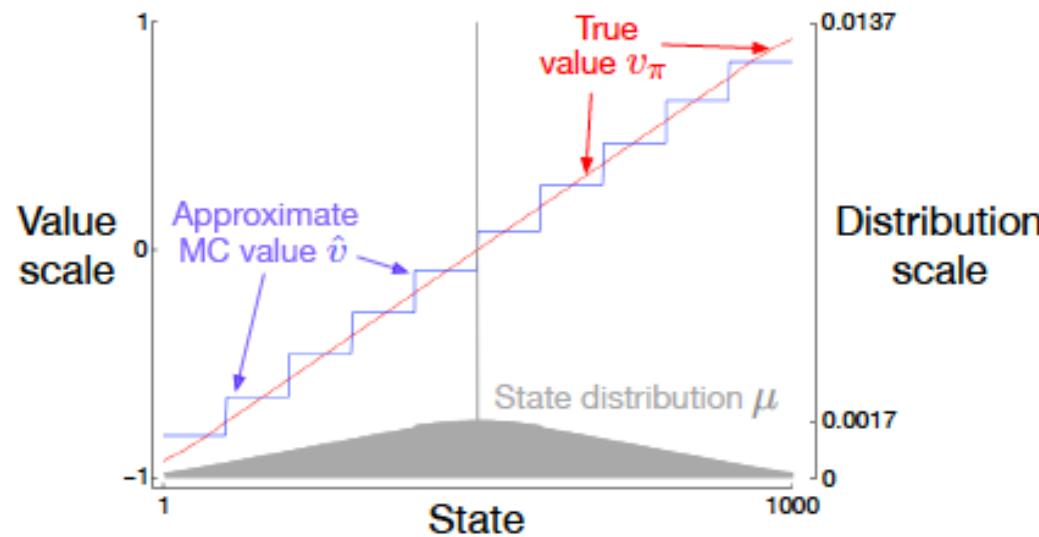
In general, less flexible function means cannot represent exactly

- Thus, we aim to approximate the true value function

An example

Let's consider just splitting the set of all states into groups

We'll only store one representative value for each group



Approximate value parametrised by this values w : $\hat{v}(s, w)$

Figure: Sutton & Barto. RL:AI

Objective

What are the ‘best’ representative values \mathbf{w} ?

$$\overline{\text{VE}}(\mathbf{w}) \doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$$



Not all states are equally important!

Objective

What are the ‘best’ representative values \mathbf{w} ?

$$\overline{VE}(\mathbf{w}) \doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$$



Not all states are equally important!

Today: On-policy distribution

- Continuing tasks: Stationary distribution

$$\mu_\pi(s) = \sum_{\bar{s}} \sum_{\mathbf{a}} p(s|\bar{s}, \mathbf{a}) \pi(\mathbf{a}|\bar{s}) \mu_\pi(\bar{s})$$

Objective

What are the ‘best’ representative values \mathbf{w} ?

$$\overline{VE}(\mathbf{w}) \doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$$



Not all states are equally important!

Today: On-policy distribution

- Continuing tasks: Stationary distribution

$$\mu_\pi(s) = \sum_{\bar{s}} \sum_a p(s|\bar{s}, a) \pi(a|\bar{s}) \mu_\pi(\bar{s})$$

- Episodic task: consider initial state distribution h :

$$\mu_\pi(s) = \frac{\eta_\pi(s)}{\sum_{s'} \eta_\pi(s')}, \quad \eta_\pi(s) = h(s) + \sum_{\bar{s}} \eta_\pi(\bar{s}) \sum_a \pi(a|\bar{s}) p(s|\bar{s}, a)$$

fraction of steps spend in s



Average number of steps spend in s per episode

Stochastic gradient descent

Let's follow the negative gradient! But how to approximate

$$\nabla_{\mathbf{w}} \overline{\text{VE}}(\mathbf{w}) = \nabla_{\mathbf{w}} \sum_{s \in \mathcal{S}} \mu(s) [v_{\pi}(s) - \hat{v}(s, \mathbf{w})]^2$$

Stochastic gradient descent

Let's follow the negative gradient! But how to approximate

$$\nabla_{\mathbf{w}} \overline{\text{VE}}(\mathbf{w}) = \nabla_{\mathbf{w}} \sum_{s \in \mathcal{S}} \mu(s) [v_{\pi}(s) - \hat{v}(s, \mathbf{w})]^2$$

- Pick random samples from $\mu(s)$ (how?)
- Generate unbiased estimates from v_{π} (how?)

Stochastic gradient descent

Let's follow the negative gradient! But how to approximate

$$\nabla_{\mathbf{w}} \overline{\text{VE}}(\mathbf{w}) = \nabla_{\mathbf{w}} \sum_{s \in \mathcal{S}} \mu(s) [v_{\pi}(s) - \hat{v}(s, \mathbf{w})]^2$$

$s_i \sim H_{\pi}(s)$

$$= \sum_i \frac{\partial}{\partial \mathbf{w}} \left[\frac{G(s_i) - \hat{V}(s_i, \mathbf{w})}{\hat{V}(s_i, \mathbf{w})} \right]^2$$

- Pick random samples from $\mu(s)$ (how?)
As μ is on-policy distribution, simply pick random visited state
- Generate unbiased estimates from v_{π} (how?)
Use the return G

Stochastic gradient descent

This leads to gradient Monte Carlo algorithm

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [G_t - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)$$

- (Use samples and returns from on-policy distribution)

Objective

In the end, we want to learn the best **policy**

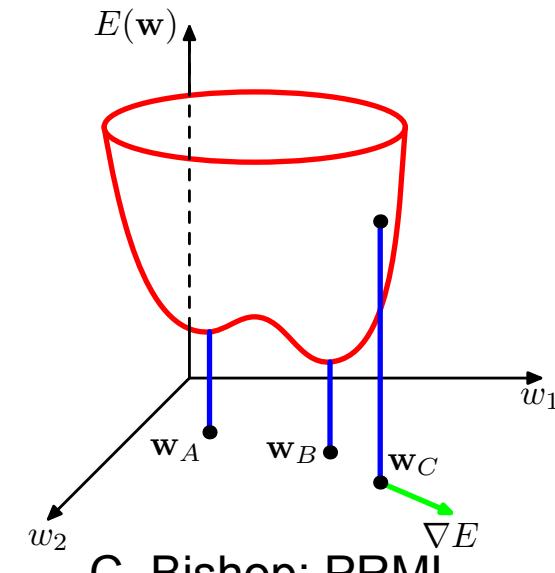
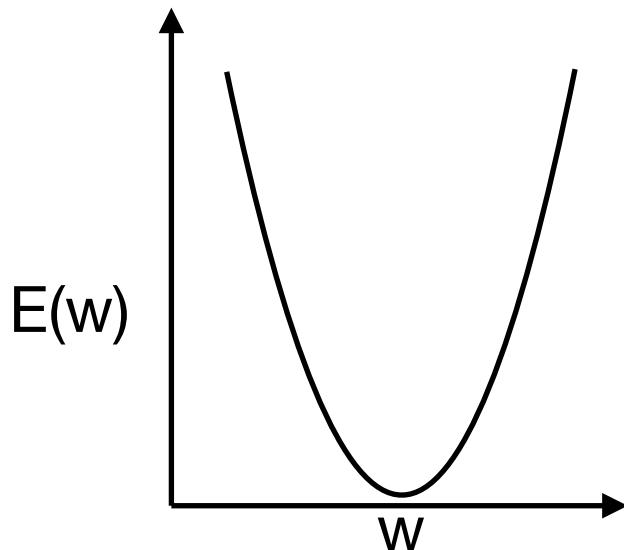
- Not clear that lowest value error leads to best policy
- But no obvious better candidate for now...

Objective

In the end, we want to learn the best **policy**

- Not clear that lowest value error leads to best policy
- But no obvious better candidate for now...

Local vs. global optima



C. Bishop; PRML

TD with function approximation

Can we also do this with bootstrapping?

Could use the bootstrapping estimate as target

$$\underline{V(s)} \leftarrow \underline{V(s)} + \alpha (\underline{V(s')} + R - \underline{V(s)})$$

$$R + \gamma \hat{v}(S', \mathbf{w})$$

Semi-gradient TD(0)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha [R + \gamma \hat{v}(S', \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)$$

TD with function approximation

Can we also do this with bootstrapping?

Could use the bootstrapping estimate as target

$$R + \gamma \hat{v}(S', \mathbf{w})$$

Semi-gradient TD(0)

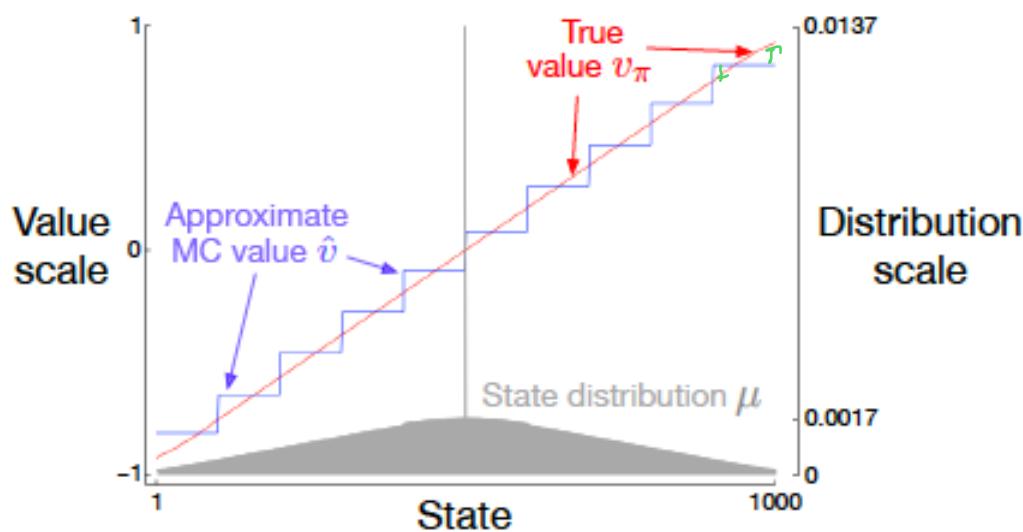
$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha [R + \gamma \hat{v}(S', \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)$$

- this ignores that the target depends on \mathbf{w} , too!
- Doesn't minimise MSTD error! Called a *semi-gradient* method
- (using the true gradient of MSTD doesn't work well - see next lecture...)

Gradient MC vs semi-gradient TD(0)

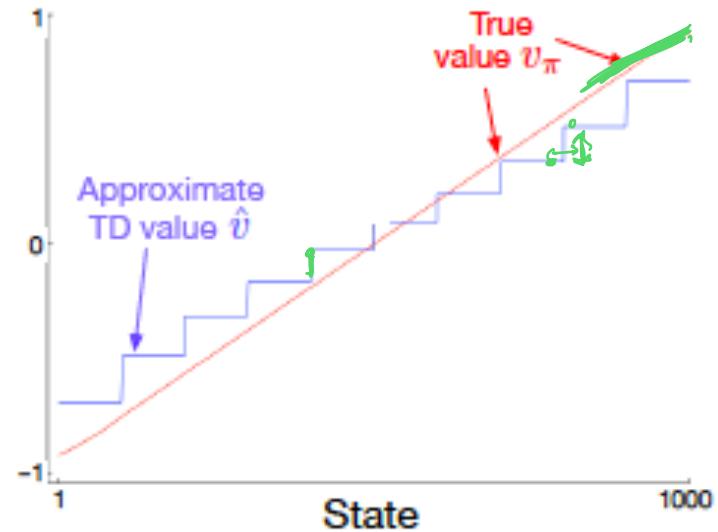
$$(\hat{V} - V_{\pi})^2$$

Gradient MC



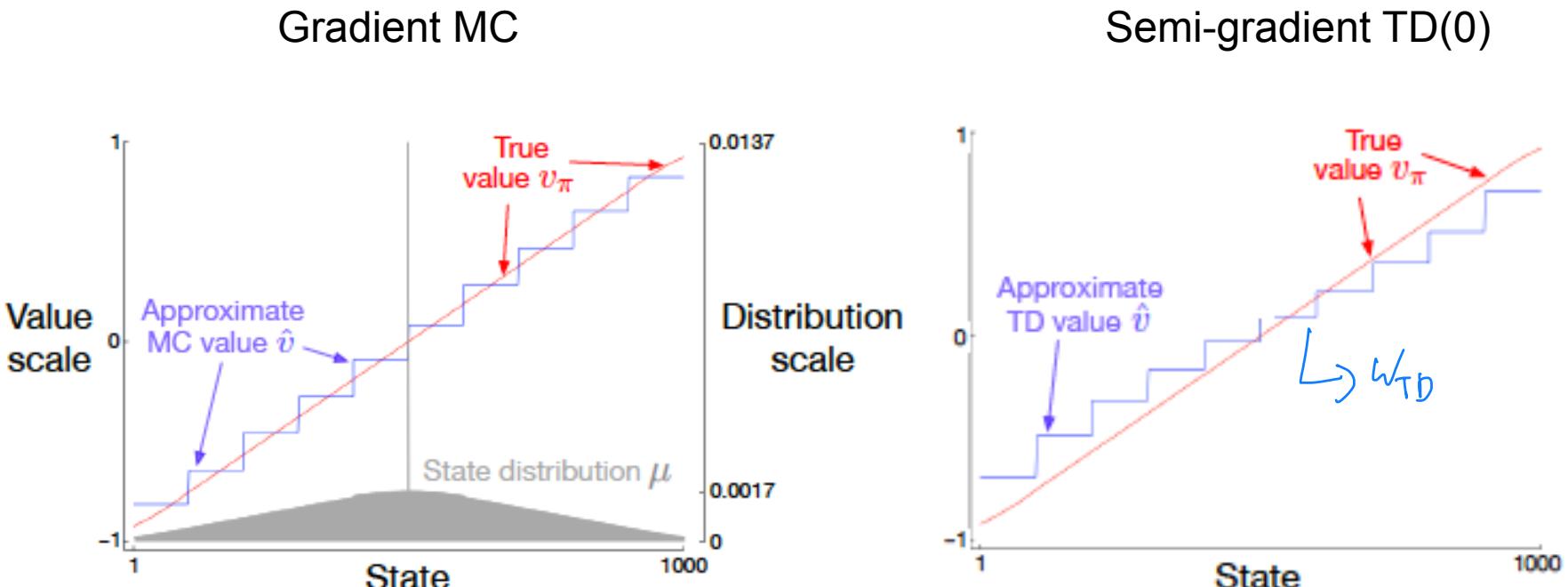
$$(V(s') - \hat{v} + R - V(s))^2$$

Semi-gradient TD(0)



Figures: Sutton & Barto. RL:AI

Gradient MC vs semi-gradient TD(0)



Figures: Sutton & Barto. RL:AI

Like before, TD has much lower variance and tends to learn faster. However, in the long run gradient MC tends to find better solutions.

Types of function approximator

$$V_{\pi} = \hat{V}(s, w)$$

Types of function approximator

Linear function approximation

- Specify some transformation $\underline{\mathbf{x}}(s) \doteq (\underline{x}_1(s), \underline{x}_2(s), \dots, \underline{x}_d(s))^T$
- Now define

$$\hat{v}(s, \mathbf{w}) \doteq \underline{\mathbf{w}}^T \underline{\mathbf{x}}(s) \doteq \sum_{i=1}^d w_i x_i(s)$$

- Note: value is linear in \mathbf{w} , but not generally in s
- We will see some useful transformation in a minute

Types of function approximator

Linear function approximation

- Specify some transformation $\mathbf{x}(s) \doteq (x_1(s), x_2(s), \dots, x_d(s))^{\top}$
- Now define

$$\hat{v}(s, \mathbf{w}) \doteq \mathbf{w}^{\top} \mathbf{x}(s) \doteq \sum_{i=1}^d w_i x_i(s)$$

- Note: value is linear in \mathbf{w} , but not generally in s
- We will see some useful transformation in a minute

Non-linear function approximation

- Any function approximation that **cannot** be brought in above form
- E.g.: (deep) neural networks
- Will be discussed later in today's lecture

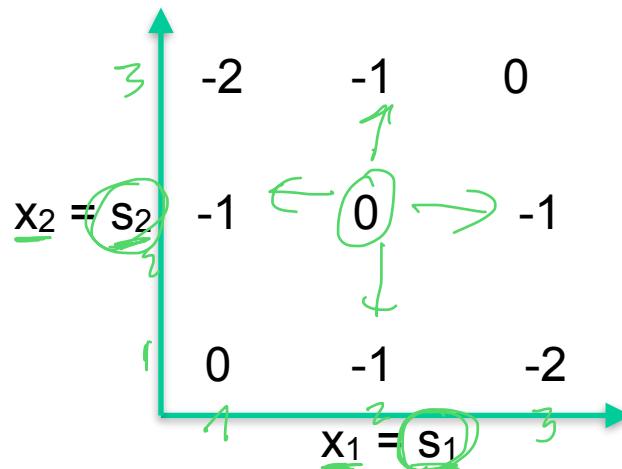
Linear function approximation

Think of \mathbf{x} as specifying important properties of the state

What is important depends on the task!

Each element of \mathbf{x} influences approximate value *linearly*, thus, no interaction between dimensions of \mathbf{x}

True value, unity transformation



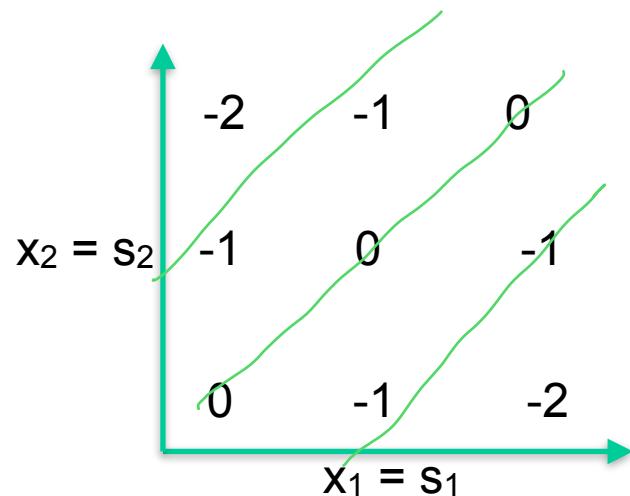
Linear function approximation

Think of \mathbf{x} as specifying important properties of the state

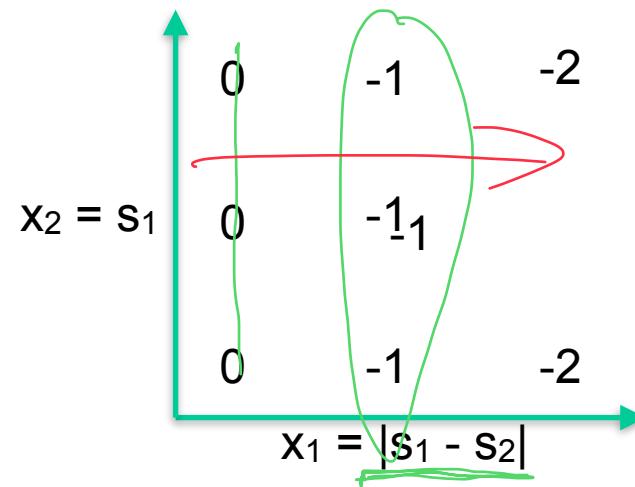
What is important depends on the task!

Each element of \mathbf{x} influences approximate value *linearly*, thus, no interaction between dimensions of \mathbf{x}

True value, unity transformation



True value, non-linear transformation



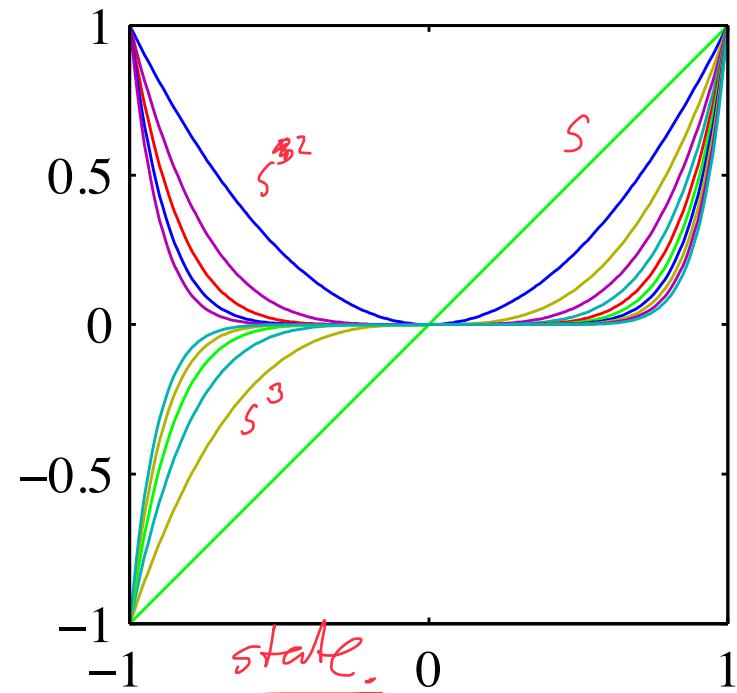
Polynomials

If we don't have task-specific transformation (*features*), use generic features that can approximate any function

Example: Polynomials

If we take “enough” polynomials,
can approximate any function

But with many features, less
generalisation



C.M. Bishop: PRML

Polynomials

Have to choose “order” of polynomial, e.g. highest power of any of the ‘raw’ state dimensions that we’ll consider

In one-d, 11th order: see image

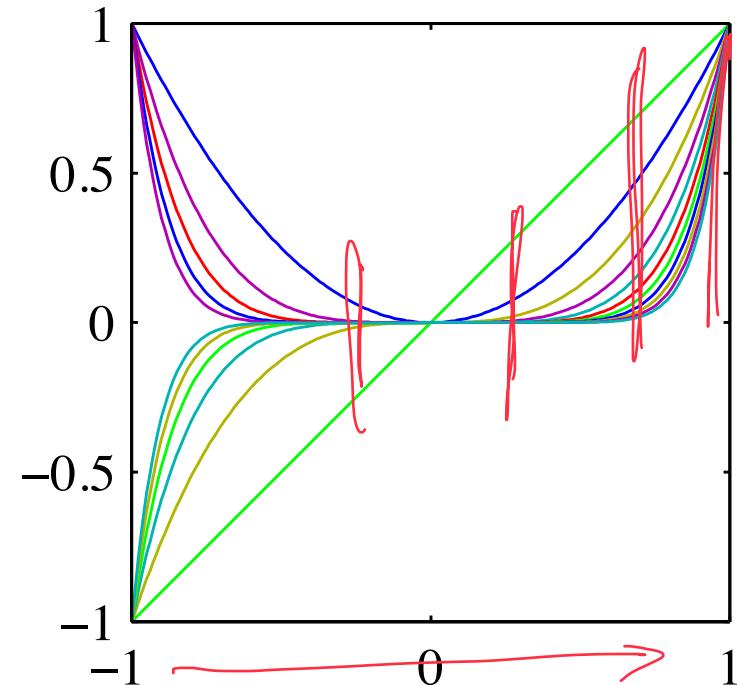
In two-d, second order

$$\mathbf{x} = [1, s_1, s_2, s_1^2, s_2^2, s_1 s_2, s_1^2 s_2, s_2^2 s_1, s_1^2 s_2^2]^T$$

Note: include ‘constant’ feature!

But:

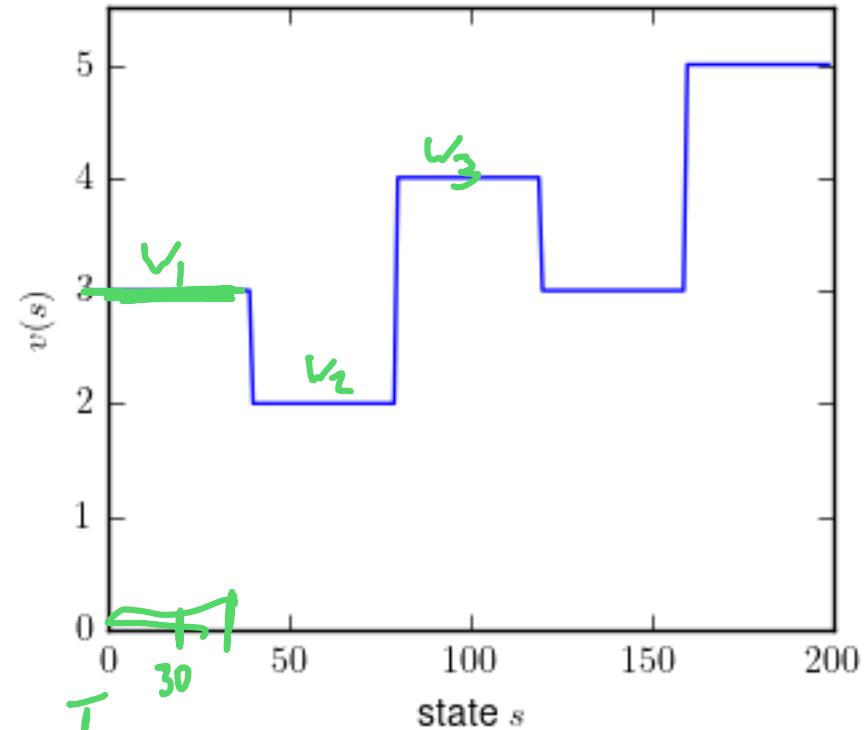
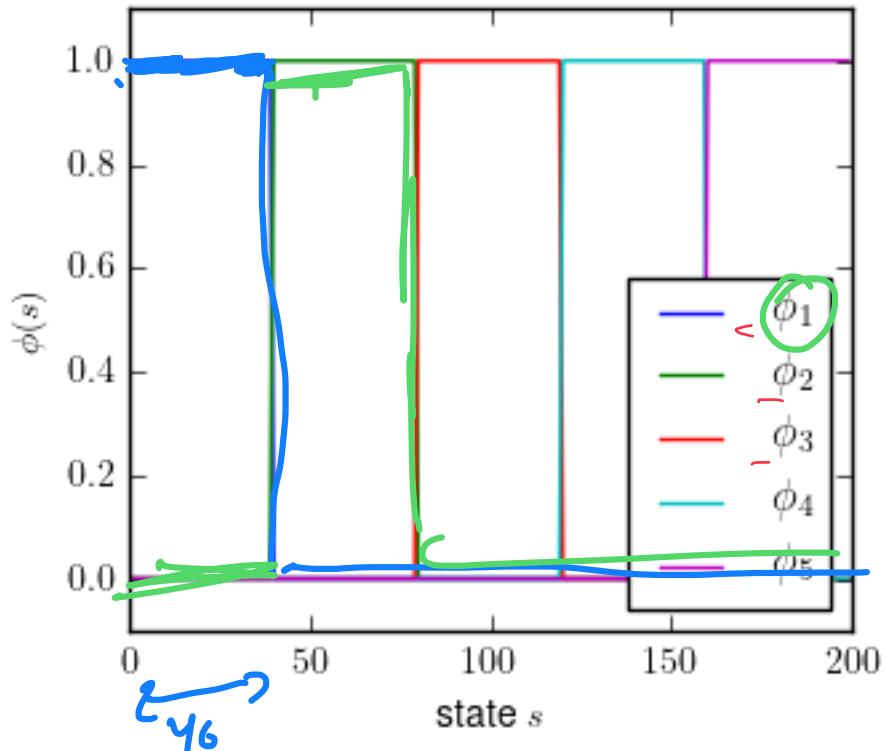
0 looks very different than extremes!



C.M. Bishop: PRML

Aggregation and tiling

$w \cdot x(s)$



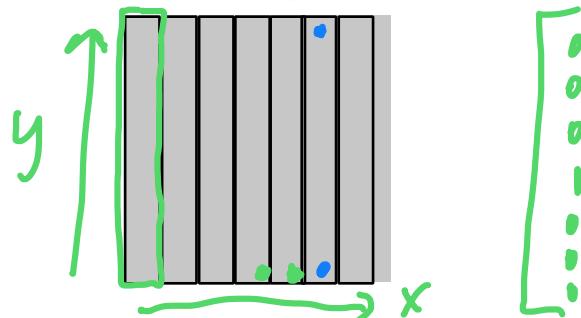
$$V(s, v) = w^T x(s) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_5 \end{bmatrix} = v_1$$

Aggregation and tiling

Different types of aggregation

x_1 : agent present in box 1 (1 or 0)

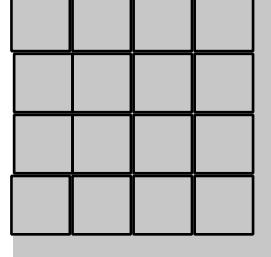
↓
larger / continuous 2-d state-space



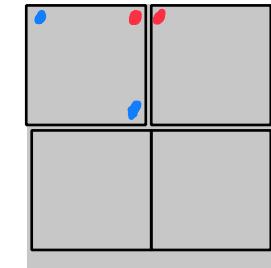
Vertical generalization: suitable
if y-coordinate does not matter

Aggregation and tiling

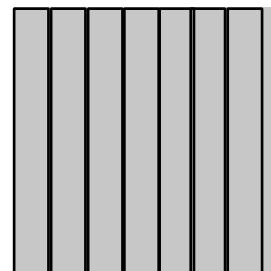
Different types of aggregation



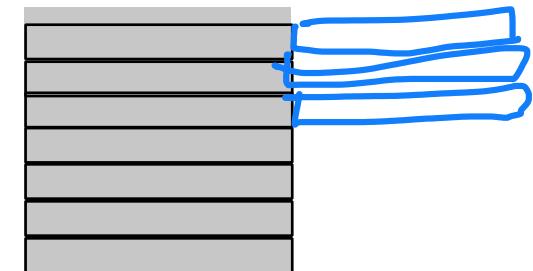
Both coordinates matter, relatively fine
discretisation (little generalization)



Both coordinates matter, relatively coarse
discretisation (high error possible)



Vertical generalization: suitable
if y-coordinate does not matter

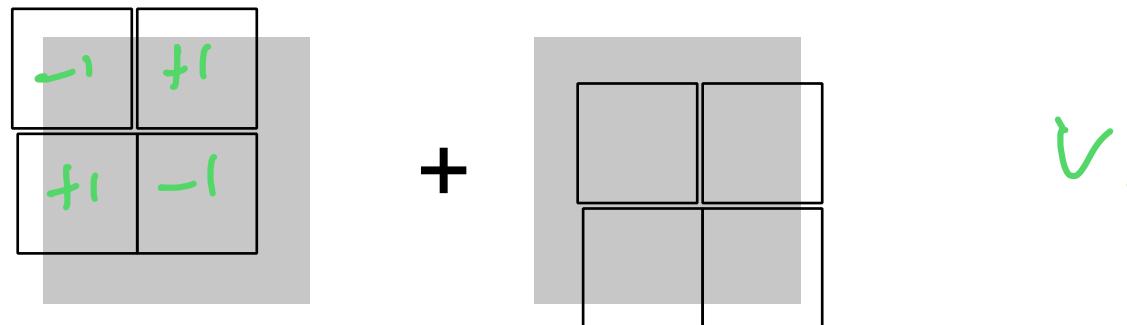


Horizontal generalization: suitable
if x-coordinate does not matter

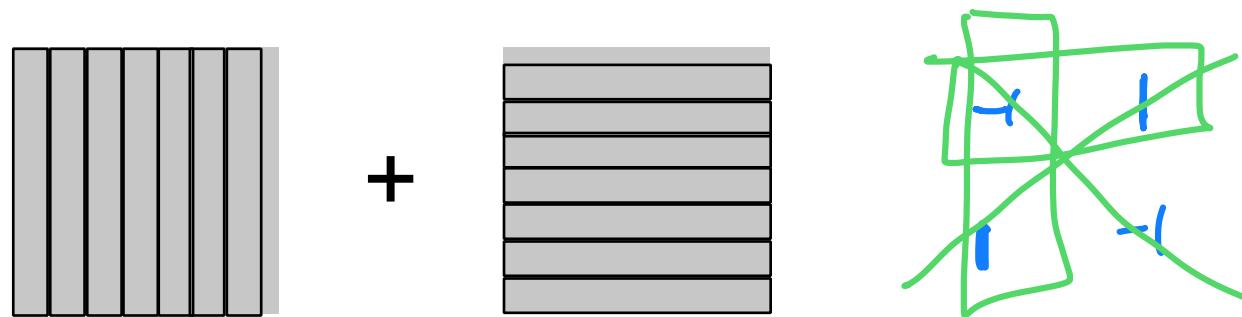
Aggregation and tiling

Pure aggregation: exactly one feature is one, all others zero

Instead: can combine multiple tiling. Exactly n features = 1



Relatively much generalization, but still able to learn functions with more detail

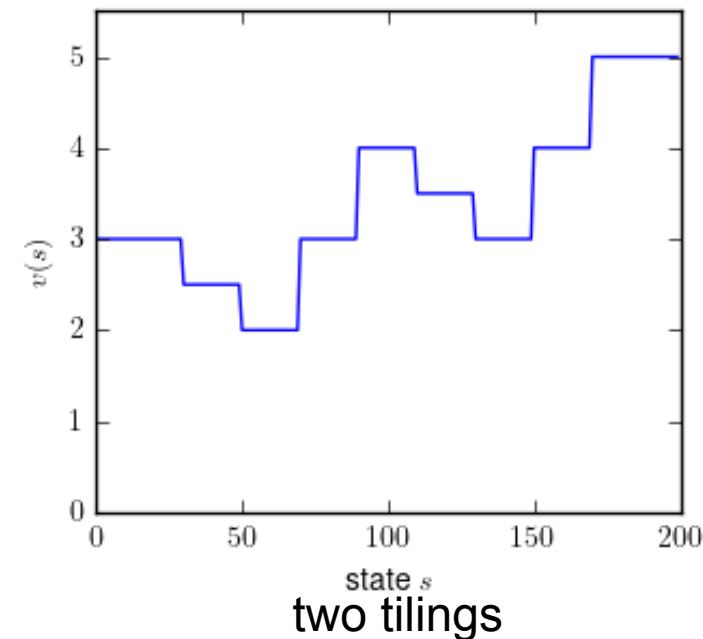
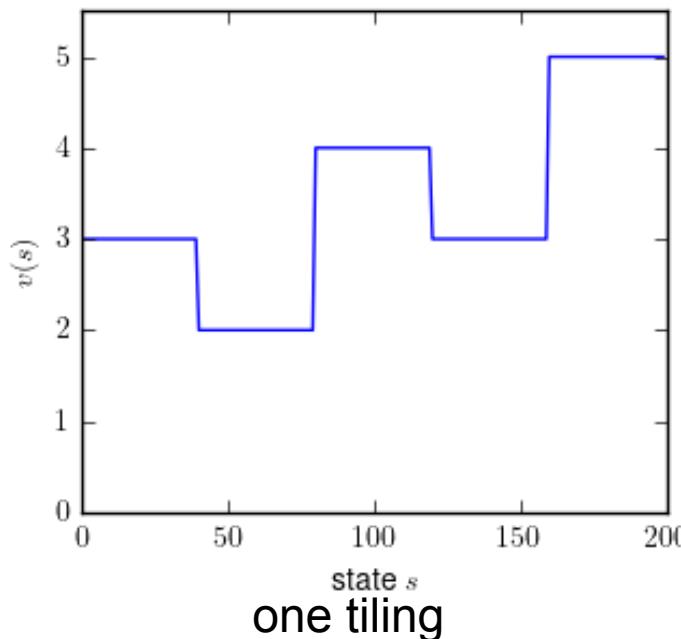


Generalize in x and y directions, but avoid modelling interaction terms

Radial basis functions

With aggregation: jump in features and value function when border in state space is crossed. Not smooth.

Tiling is a bit better. We now make small jumps for each of several aggregations.



Radial basis functions

With aggregation: jump in features and value function when border in state space is crossed. Not smooth.

Tiling is a bit better. We now make small jumps for each of several aggregations.

But: we can have features where we go from 0 to 1 smoothly

Radial basis features:

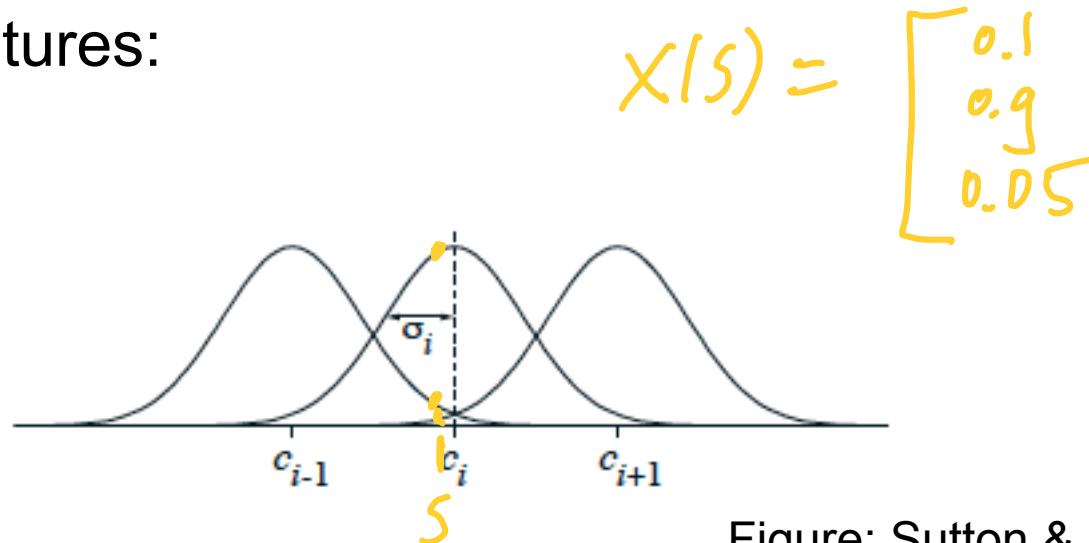


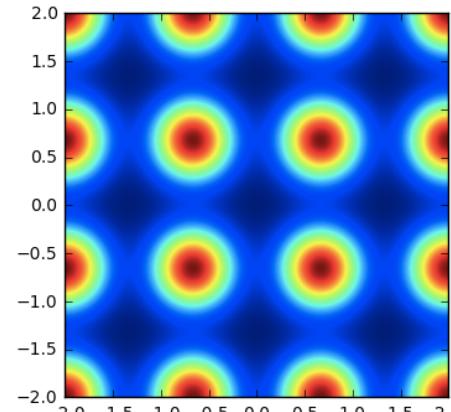
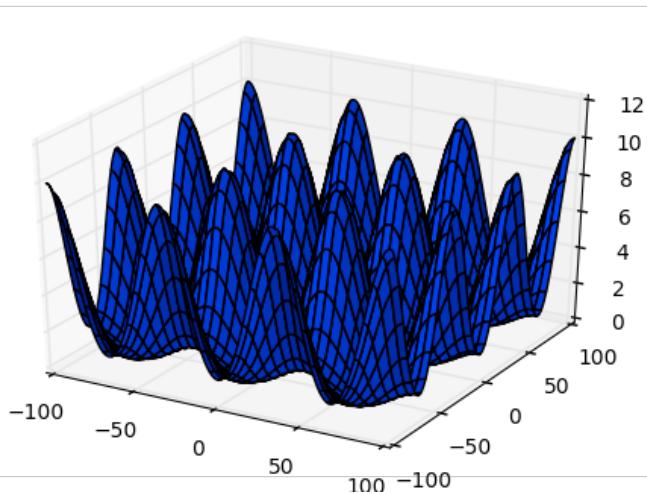
Figure: Sutton & Barto. RL:AI

Radial basis functions

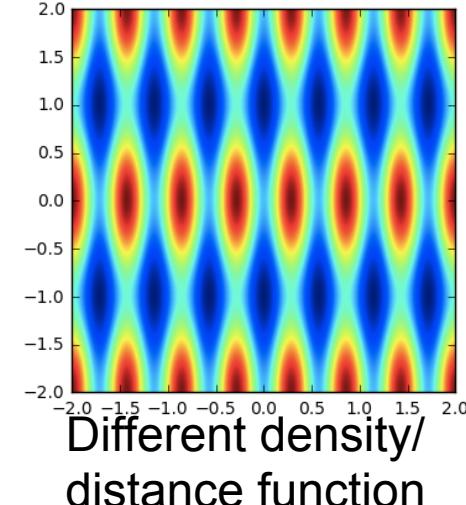
Radial basis functions only depend on distance to center, e.g.:

$$x_i(s) \doteq \exp\left(-\frac{\|s - c_i\|^2}{2\sigma_i^2}\right)$$

Also work in multi-d



Top-down view

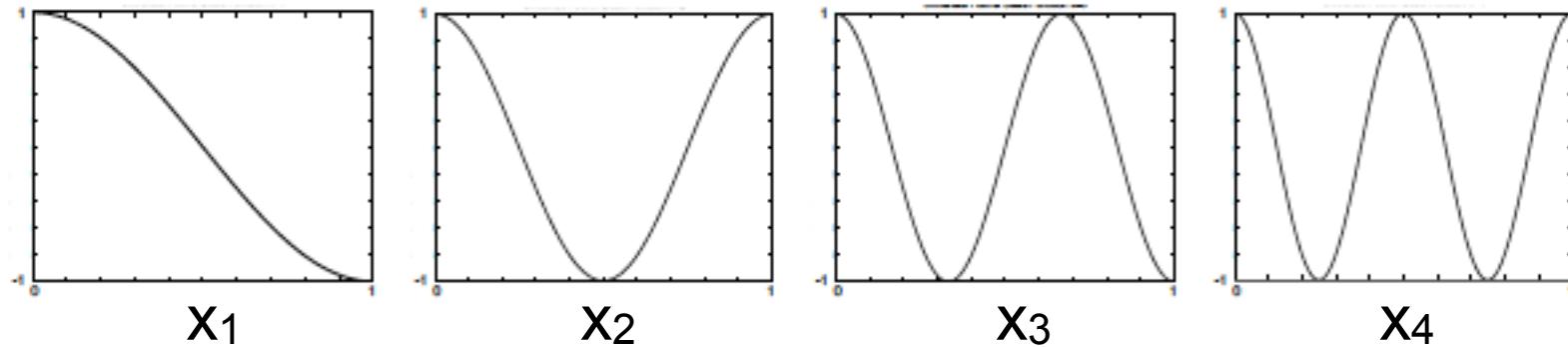


Different density/
distance function

Fourier basis

First: scale raw states between 0 and 1

Then (one-d): $x_j(s) = \cos(j\pi s)$, $s \in [0, 1]$
(so x_0 again a straight line!)



Again, pick as many features as needed...

Figure: Sutton & Barto. RL:AI

Linear function approximation continued

Any of these choices leads to linear approximations to the value function, i.e.:

$$\hat{v}(s, \mathbf{w}) \doteq \mathbf{w}^\top \mathbf{x}(s) \doteq \sum_{i=1}^d w_i x_i(s)$$

for different choices of the *features* \mathbf{x}

Convince yourself:

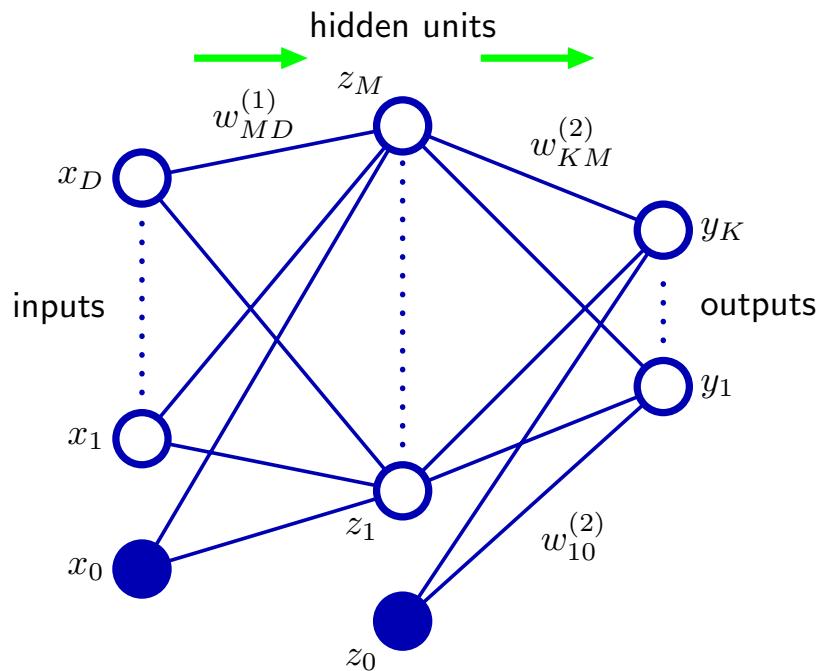
- What does \mathbf{x} look like for the tilings example?
- Is tabular RL a specific case of a linear function? What would \mathbf{x} be?

Non-linear function approximation

Any other type of function approximation in non-linear

A popular non-linear approximator is a neural network

E.g. feedforward network:



$$\hat{v}(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}) = \sum_{m=0}^M w_m^{(2)} h^{(1)} \left(\sum_{d=0}^D w_{md}^{(1)} x_d \right)$$

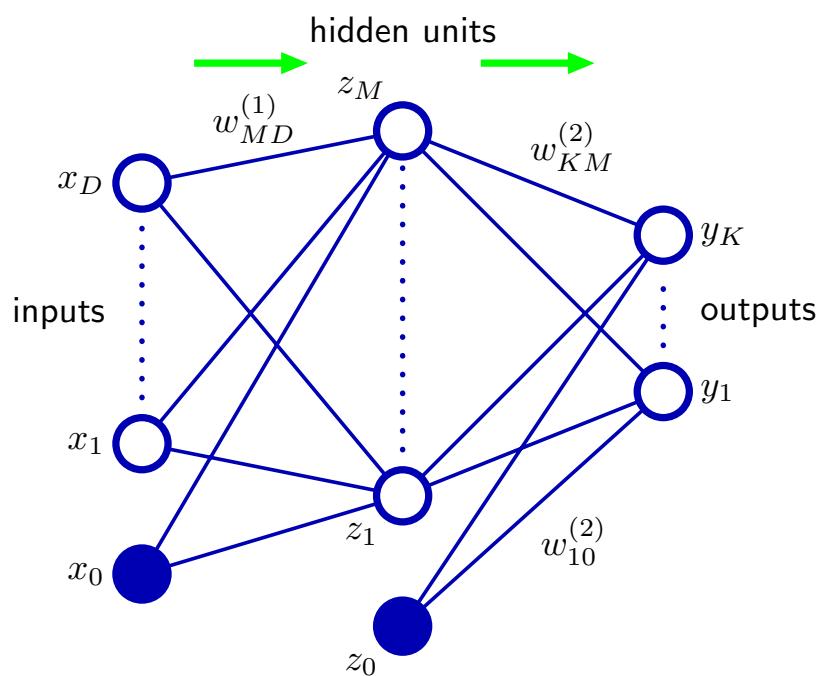
Non-linear activation function

Non-linear function approximation

Any other type of function approximation in non-linear

A popular non-linear approximator is a neural network

E.g. feedforward network:



‘weights’ \downarrow ‘features’ \uparrow

$$\hat{v} \left(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)} \right) = \sum_{m=0}^M w_m^{(2)} h^{(1)}$$

$\left(\sum_{d=0}^D w_{md}^{(1)} x_d \right)$

Non-linear activation function

Using the defined features

Recall the definition of gradient Monte Carlo

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha [G_t - \hat{v}(S_t, \mathbf{w}_t)] \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w}_t)$$

and semi-gradient TD(0)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha [R + \gamma \hat{v}(S', \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w}_t)$$

In both cases, we can simply ‘plug in’ known quantities and estimated values, but we still need $\nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$

Using the defined features

For linear function approximation extremely easy!

$$\nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w}) = \nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{x}(s) = \mathbf{x}(s)$$

For non-linear approximation a bit more involved

- E.g. neural network

$$\hat{v}(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}) = \sum_{m=0}^M w_m^{(2)} h^{(1)} \left(\underbrace{\sum_{d=0}^D w_{md}^{(1)} x_d}_{a_m} \right)$$

- Can figure this out: repeated application of chain rule, etc (backpropagation)
- Easy with auto-diff frameworks (PyTorch, TensorFlow, ...)

Using the defined features

$$g \delta V(S_{t+1}) + R$$

Linear function approximation very nice to work with

- Easy to calculate required gradient
 - Update convenient form $\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [U_t - \hat{v}(S_t, \mathbf{w}_t)] \mathbf{x}(S_t)$
 - Can prove: all local optima are global optima
 - If features linearly independent: only one optimum
-
- So: gradient MC always converges to global minimum of VE with linear approximation
 - With non-linear functions, there might be local minima
 - What about semi-gradient TD(0)?

Fixed point of semi-gradient TD(0)

$$\begin{aligned} \mathbf{w}_{t+1} &\doteq \mathbf{w}_t + \alpha \left(R_{t+1} + \gamma \mathbf{x}_{t+1}^\top \mathbf{w}_t - \mathbf{x}_t^\top \mathbf{w}_t \right) \mathbf{x}_t \\ &= \mathbf{w}_t + \alpha \left(R_{t+1} \mathbf{x}_t - \mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top \mathbf{w}_t \right) \end{aligned}$$

$$E[\mathbf{w}_{t+1} | \mathbf{w}_t] = \mathbf{w}_t + \alpha \left(E[R_{t+1} \mathbf{x}_t] - E[\mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top] \mathbf{w}_t \right)$$

$\underbrace{\hspace{10em}}_b$ $\underbrace{\hspace{10em}}_A$

$$\mathbf{w}_{t+1} = \mathbf{w}_t$$

$$b - A \mathbf{w}_{TD} = 0$$

$$b \doteq A \mathbf{w}_{TD}$$

$$\mathbf{w}_{TD} \doteq \underline{A^{-1} b}$$

Fixed point of semi-gradient TD(0)

$$\begin{aligned}\mathbf{w}_{t+1} &\doteq \mathbf{w}_t + \alpha \left(R_{t+1} + \gamma \mathbf{x}_{t+1}^\top \mathbf{w}_t - \mathbf{x}_t^\top \mathbf{w}_t \right) \mathbf{x}_t \\ &= \mathbf{w}_t + \alpha \left(R_{t+1} \mathbf{x}_t - \mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top \mathbf{w}_t \right)\end{aligned}$$

$$\mathbb{E}[\mathbf{w}_{t+1} | \mathbf{w}_t] = \mathbf{w}_t + \alpha \left(\underbrace{\mathbb{E}[R_{t+1} \mathbf{x}_t]}_{\mathbf{b}} - \underbrace{\mathbb{E}[\mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top]}_{\mathbf{A}} \mathbf{w}_t \right)$$

At convergence $\mathbf{w}_{t+1} = \mathbf{w}_t$, so

$$\mathbf{b} - \mathbf{A}\mathbf{w}_{\text{TD}} = \mathbf{0}$$

$$\mathbf{b} = \mathbf{A}\mathbf{w}_{\text{TD}}$$

$$\mathbf{w}_{\text{TD}} \doteq \mathbf{A}^{-1}\mathbf{b}, \quad \text{the TD fixed point}$$

Fixed point of semi-gradient TD(0)

Semi-gradient TD with linear features can be proven to **converge** to this fixed point w_{+b}
(when on-policy, independent features)

With non-linear features, semi-gradient TD might **diverge**

Fixed point of semi-gradient TD(0)

Semi-gradient TD with linear features can be proven to **converge** to this fixed point
(when on-policy, independent features)

With non-linear features, semi-gradient TD might **diverge**

Unfortunately, w_{TD} is generally not the minimiser of VE

What can be proven (under some conditions) is that

$$\overline{VE}(w_{TD}) \leq \frac{1}{1-\gamma} \min_w \overline{VE}(w)$$

= solution Gradient MC.

“maximally $1/(1-\gamma)$ worse than the best possible error”

- not very good for γ close to 1!

Fixed point of semi-gradient TD(0)

Semi-gradient TD with linear features can be proven to **converge** to this fixed point
(when on-policy, independent features)

With non-linear features, semi-gradient TD might **diverge**

Unfortunately, w_{TD} is generally not the minimiser of VE

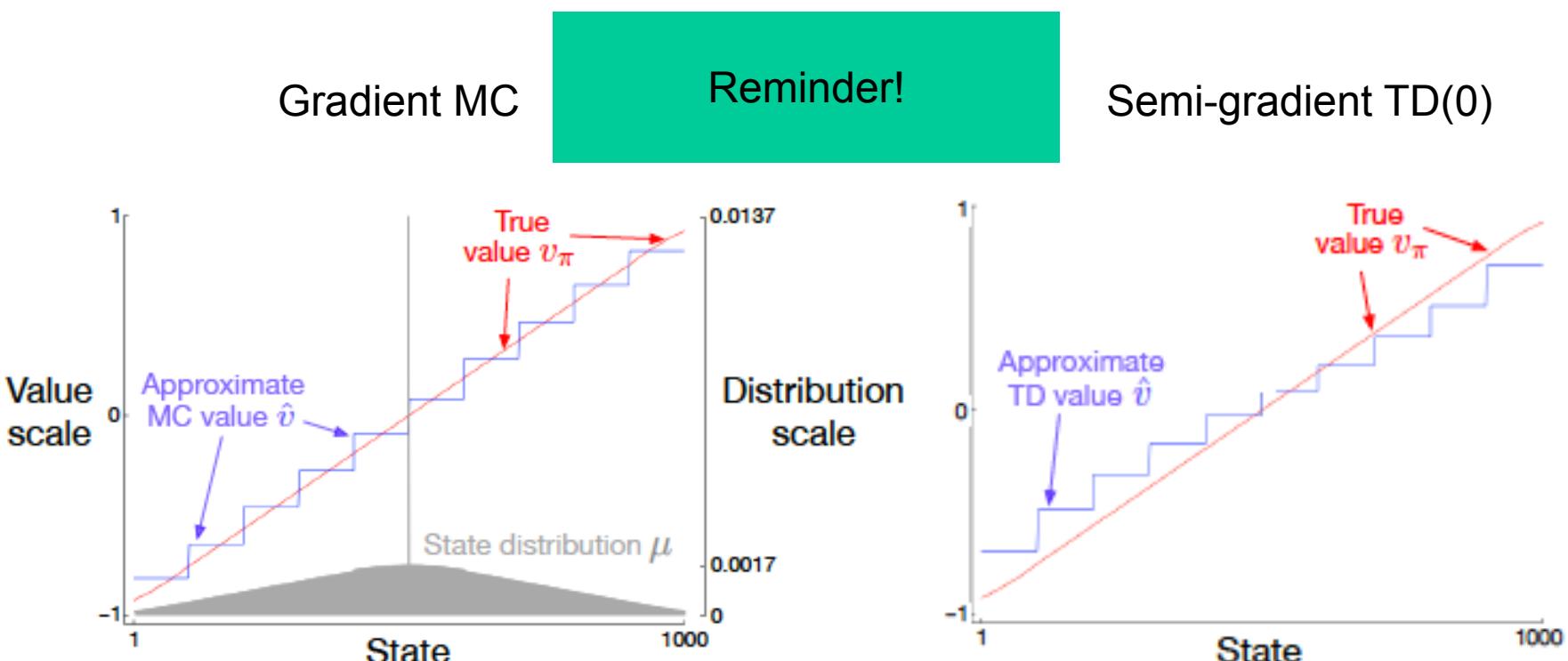
What can be proven (under some conditions) is that

$$\overline{VE}(w_{TD}) \leq \frac{1}{1-\gamma} \min_w \overline{VE}(w)$$

“maximally $1/(1-\gamma)$ worse than the best possible error”
- not very good for γ close to 1!

Still, semi-gradient TD converges much faster. Within a fixed amount of samples, TD can still beat MC...

Gradient MC vs semi-gradient TD(0)



Figures: Sutton & Barto. RL:AI

Like before, TD has much lower variance and tends to learn faster. However, if we can train however long we want the final results can be worse.

Selecting step size

Selecting step size is hard!

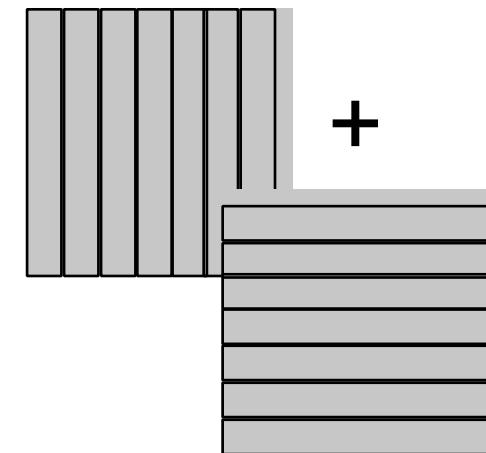
→ Tabular intuition: $\alpha=1$ jumps to best solution for this experience

To average over, say, 10 experiences use $\alpha \approx 0.1$

Function approximation bit more tricky. With e.g. two tilings and $\alpha=1$, where do we end up?

(reminder:)

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [U_t - \hat{v}(S_t, \mathbf{w}_t)] \mathbf{x}(S_t)$$



Selecting step size

Take the typical norm of gradient into account?

Rule of thumb:

$$\alpha \doteq \left(\tau \mathbb{E} \underbrace{[\mathbf{x}^\top \mathbf{x}]}_{\mathcal{Z}} \right)^{-1}$$

0.1
0.05

where τ is the number of experiences to average over.

Works best if the norm of features is close to constant

- for which discussed features is that the case?
- what could we do for the other types of features?

Least-squares temporal-difference (LSTD)

Can we get around defining stepsizes?

Try to directly find the point gradient descent will converge to...

Least-squares temporal-difference (LSTD)

Can we get around defining stepsizes?

Try to directly find the point gradient descent will converge to...

Recall the TD fixed point

$$\mathbb{E}[\mathbf{w}_{t+1} | \mathbf{w}_t] = \mathbf{w}_t + \alpha \left(\underbrace{\mathbb{E}[R_{t+1} \mathbf{x}_t]}_{\mathbf{b}} - \underbrace{\mathbb{E}[\mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top]}_{\mathbf{A}} \mathbf{w}_t \right)$$

$$\mathbf{w}_{\text{TD}} \doteq \mathbf{A}^{-1} \mathbf{b},$$

doesn't depend on step size!

Try to estimate **b** and **A** *directly from data*?

Least-squares temporal-difference (LSTD)

Use sample averages/sums to estimate

$$\mathbf{b} = \mathbb{E}[\underline{R_{t+1}\mathbf{x}_t}], \quad \mathbf{A} = \mathbb{E}[\underline{\mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top}]$$

then (ignoring the $1/(t-1)$ term):

$$\hat{\mathbf{A}}_t \doteq \sum_{k=0}^{t-1} \mathbf{x}_k (\mathbf{x}_k - \gamma \mathbf{x}_{k+1})^\top + \underbrace{\varepsilon \mathbf{I}}_{\leftarrow} \quad \text{and} \quad \hat{\mathbf{b}}_t \doteq \sum_{k=0}^{t-1} \underline{R_{k+1}\mathbf{x}_k}$$

then

$$\underline{\mathbf{w}_t} \doteq \underline{\hat{\mathbf{A}}_t^{-1} \hat{\mathbf{b}}_t}$$

Least-squares temporal-difference (LSTD)

More sample efficient than semi-gradient TD

More computationally intensive

- Naive inversion is cubic in d , but we can do it incrementally (quadratic)
- Semi-gradient was $O(d)$ in memory and computation

Doesn't require step size parameter

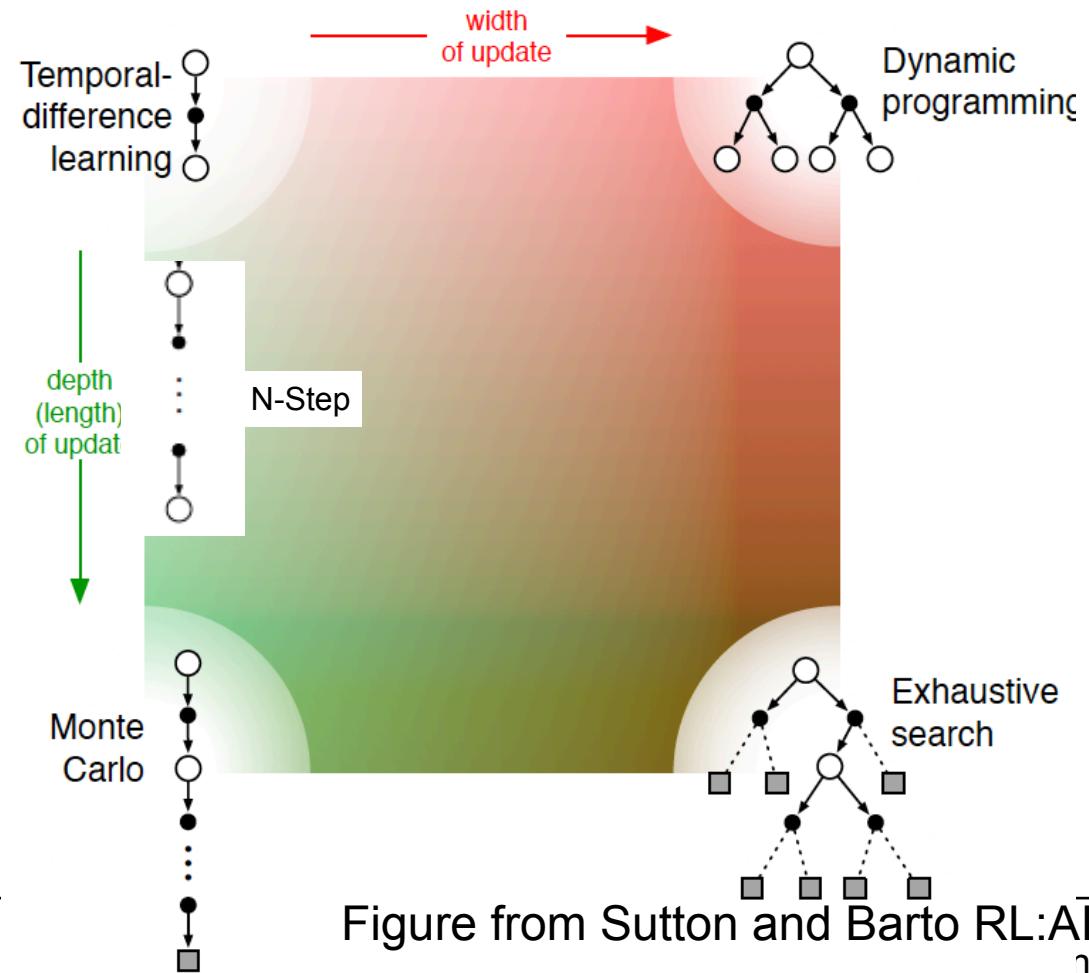
- But it does require setting ϵ

LSTD never forgets

- Can be good, but means system cannot change!

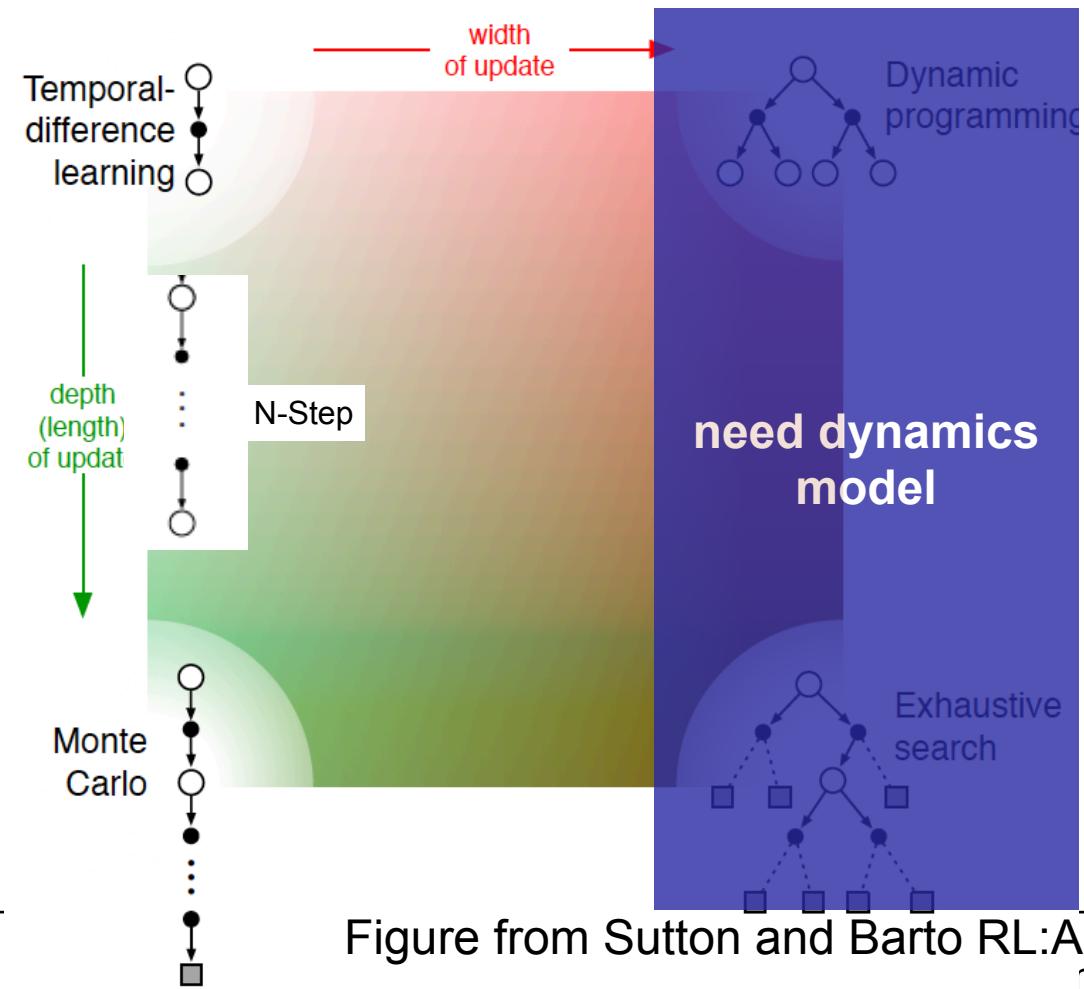
Back to the big picture

We have talked today about *prediction* methods



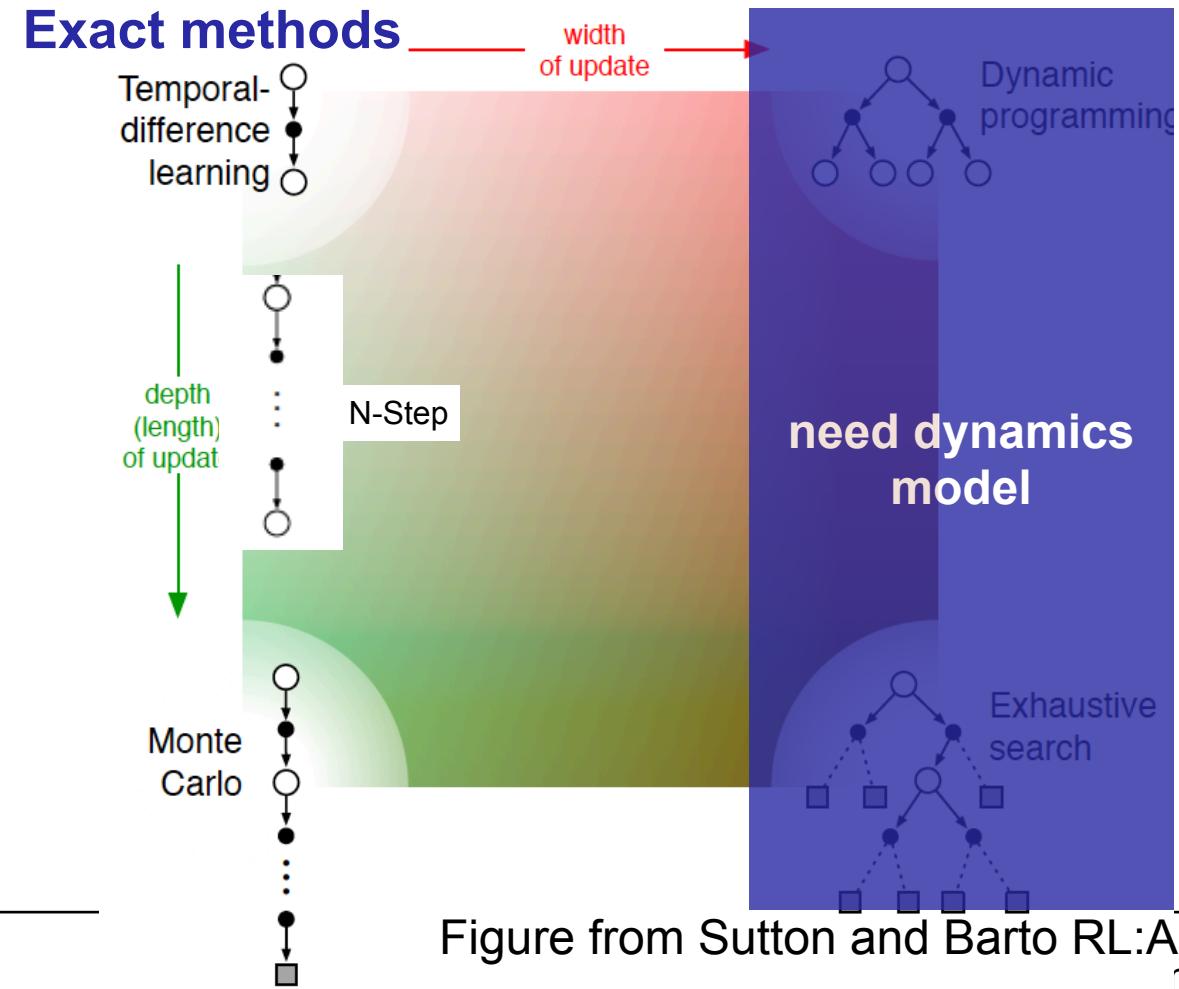
Back to the big picture

We have talked today about *prediction* methods



Back to the big picture

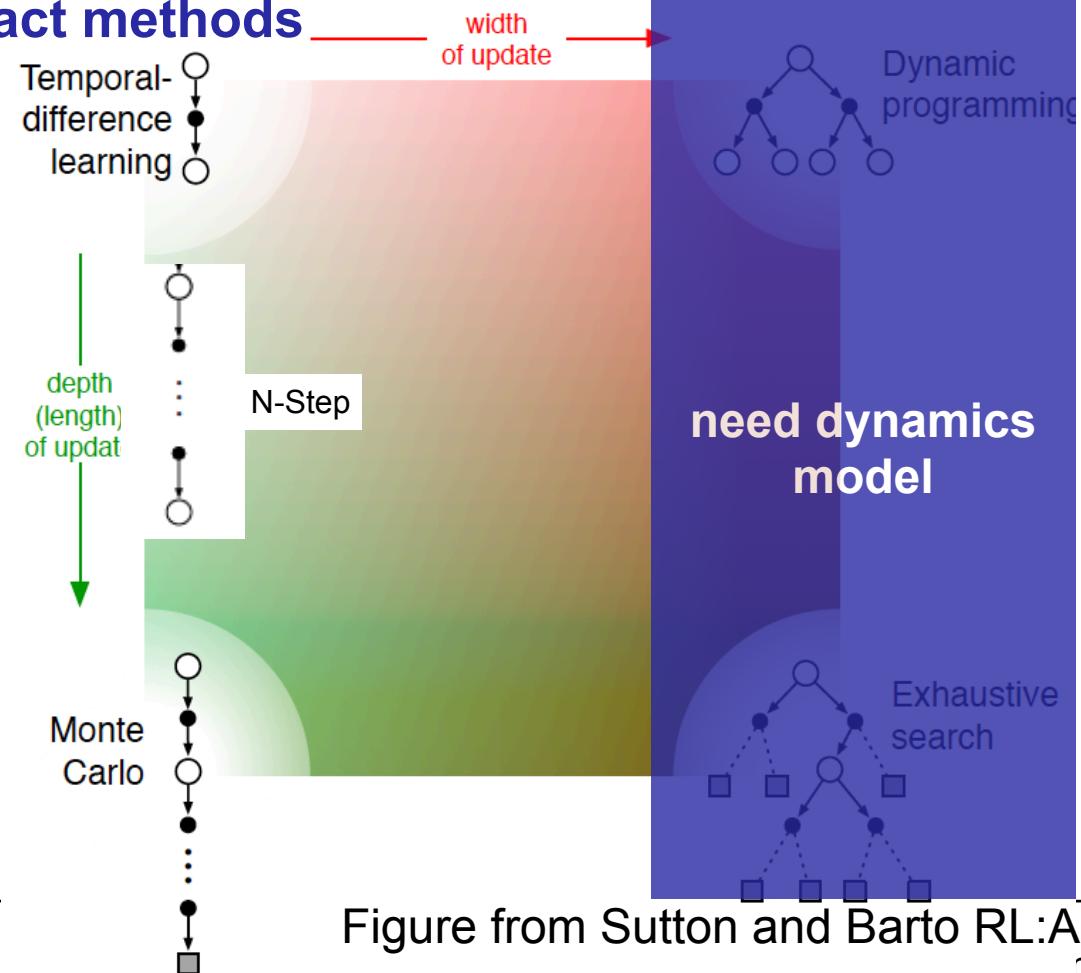
We have talked today about *prediction* methods



Back to the big picture

We have talked today about *prediction* methods

Approximate method Exact methods



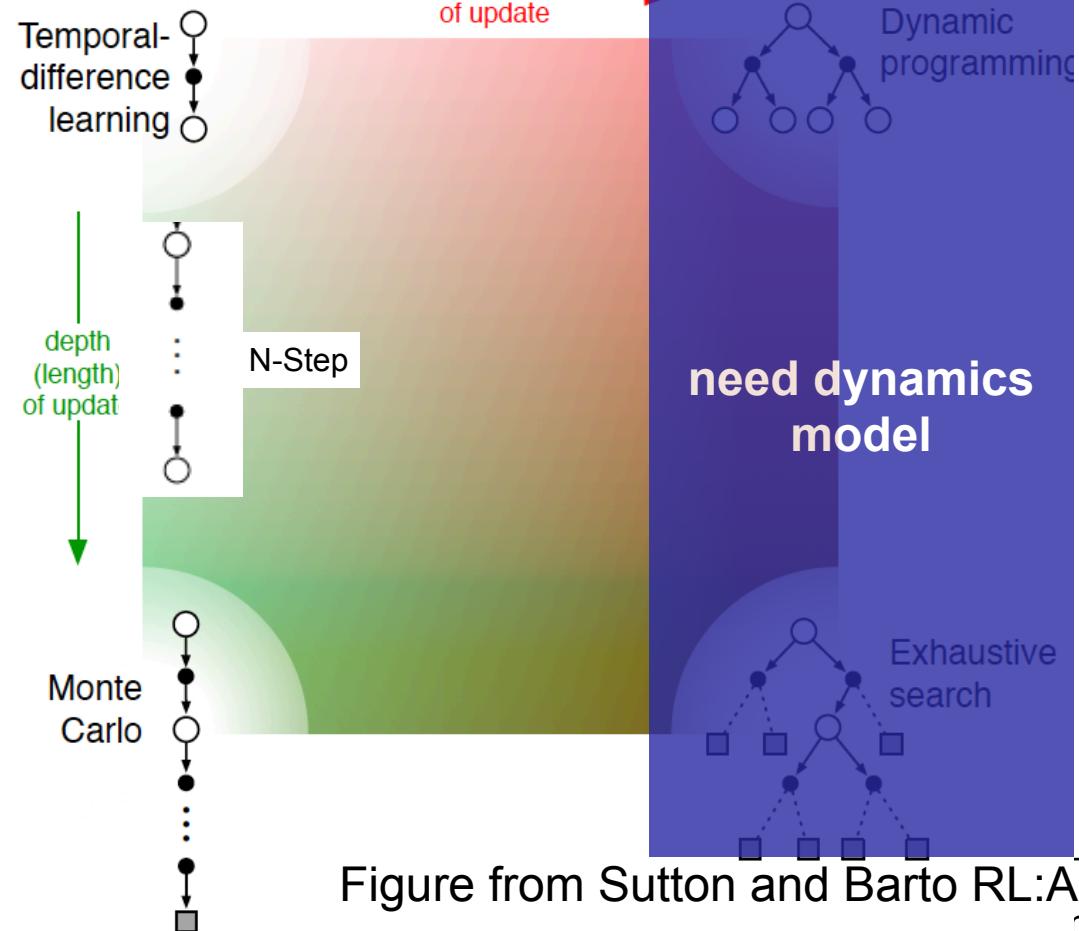
Back to the big picture

We have talked today about *prediction* methods

Approximate method

Exact methods

gradient
Monte-Carlo



Back to the big picture

We have talked today about *prediction* methods

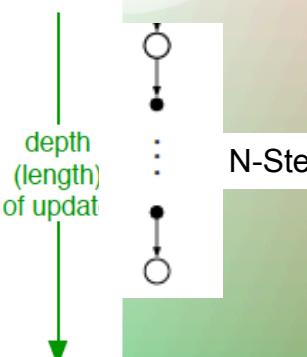
Approximate method

semi-gradient
TD(0)
& LSTD

Exact methods

Temporal-difference learning

width of update

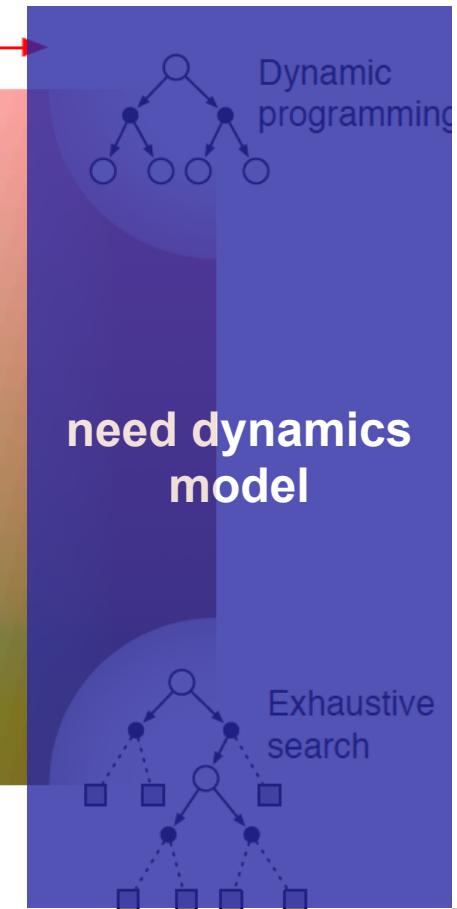


gradient
Monte-Carlo

Monte Carlo

...
↓

Figure from Sutton and Barto RL:AI
ng



Back to the big picture

We have talked today about *prediction* methods

Approximate method

semi-gradient
TD(0)
& LSTD

↓
Straightforward
generalisation
to n-step
(not discussed)

gradient
Monte-Carlo

Exact methods

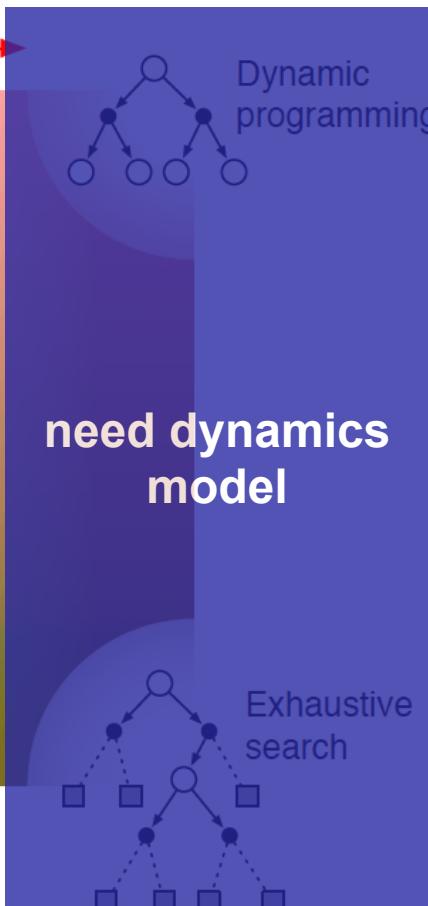
Temporal-
difference
learning

N-Step

Monte
Carlo

width
of update

depth
(length)
of update



Back to the big picture

Guarantees depend on

target computed from MC is = G_t is unbiased
target computed from TD (semi gradient TD or LSTD)
is biased.

method

	Tabular & Linear function approximator **	Nonlinear function approximator
Gradient MC	Convergence* to a <i>global minimum of VE</i>	Convergence* to a <i>local minimum of VE</i>
Semi-gradient TD	Convergence* to TD fixed point	No guarantee of convergence
LSTD	Convergence to TD fixed point	-

* with on-policy data and appropriate step-size schedule

** if features independent, single solution

Back to the big picture

For any of the methods (gradient MC / semi-gradient TD),
choice of function approximation

linear

non-linear

tabular
aggregate
tiling
radial basis function
polynomial basis function
fourier basis function

e.g. neural network

Back to the big picture

Main problems:

- prediction
- on-policy control
- off-policy control

Back to the big picture

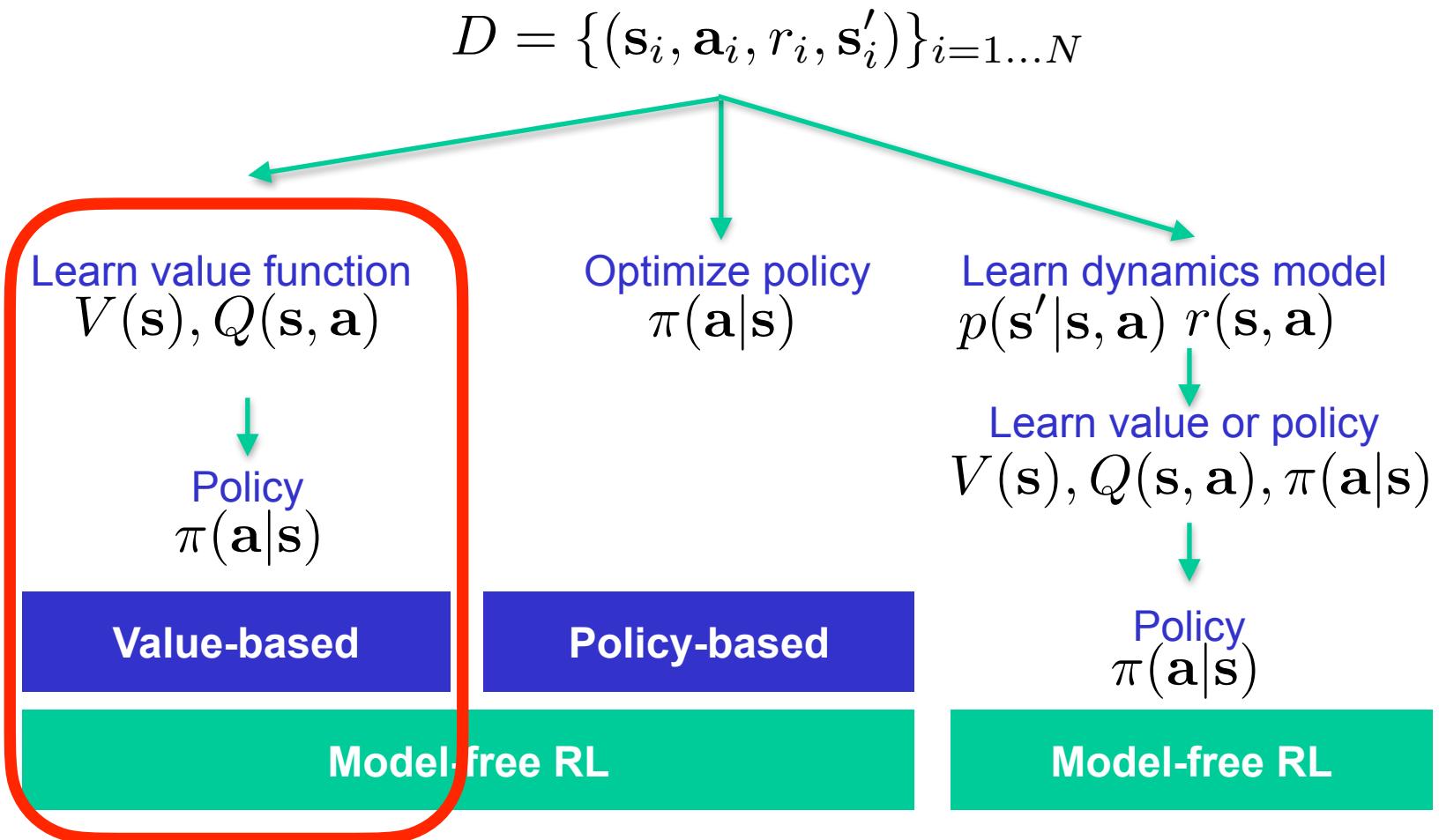
Main problems:

- prediction
- on-policy control
- off-policy control

Tabular (Exact)
lecture 2, 3& 4
lecture 2, 3& 4
lecture 2, 3& 4

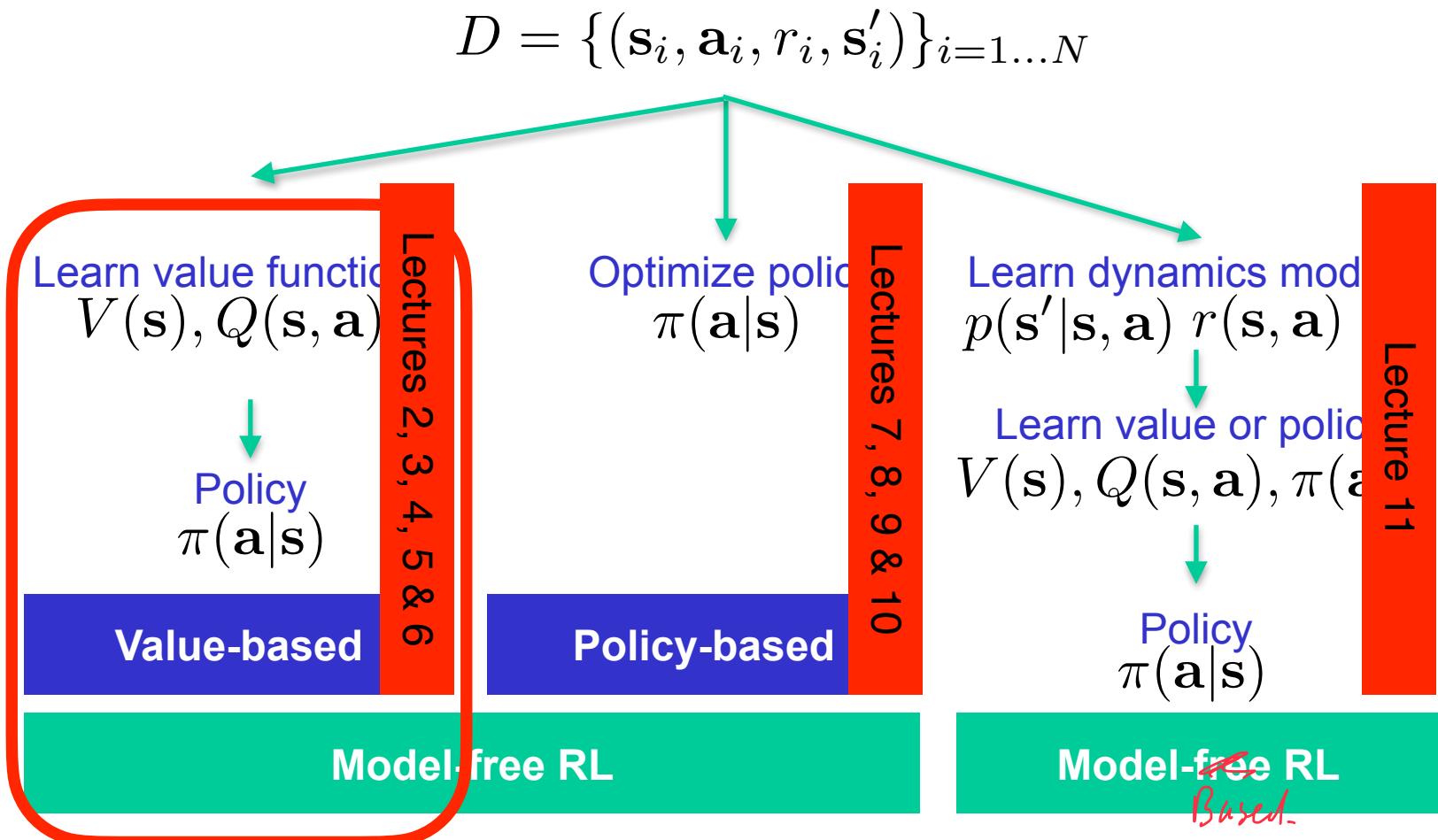
Approximate
Today's lecture (5)
Next lecture (6)
Next lecture (6)

Back to the big picture



Thanks to Jan Peters

Back to the big picture



Thanks to Jan Peters

Feedback?

h.c.vanhoof@uva.nl

Extra: Aggregation and tiling

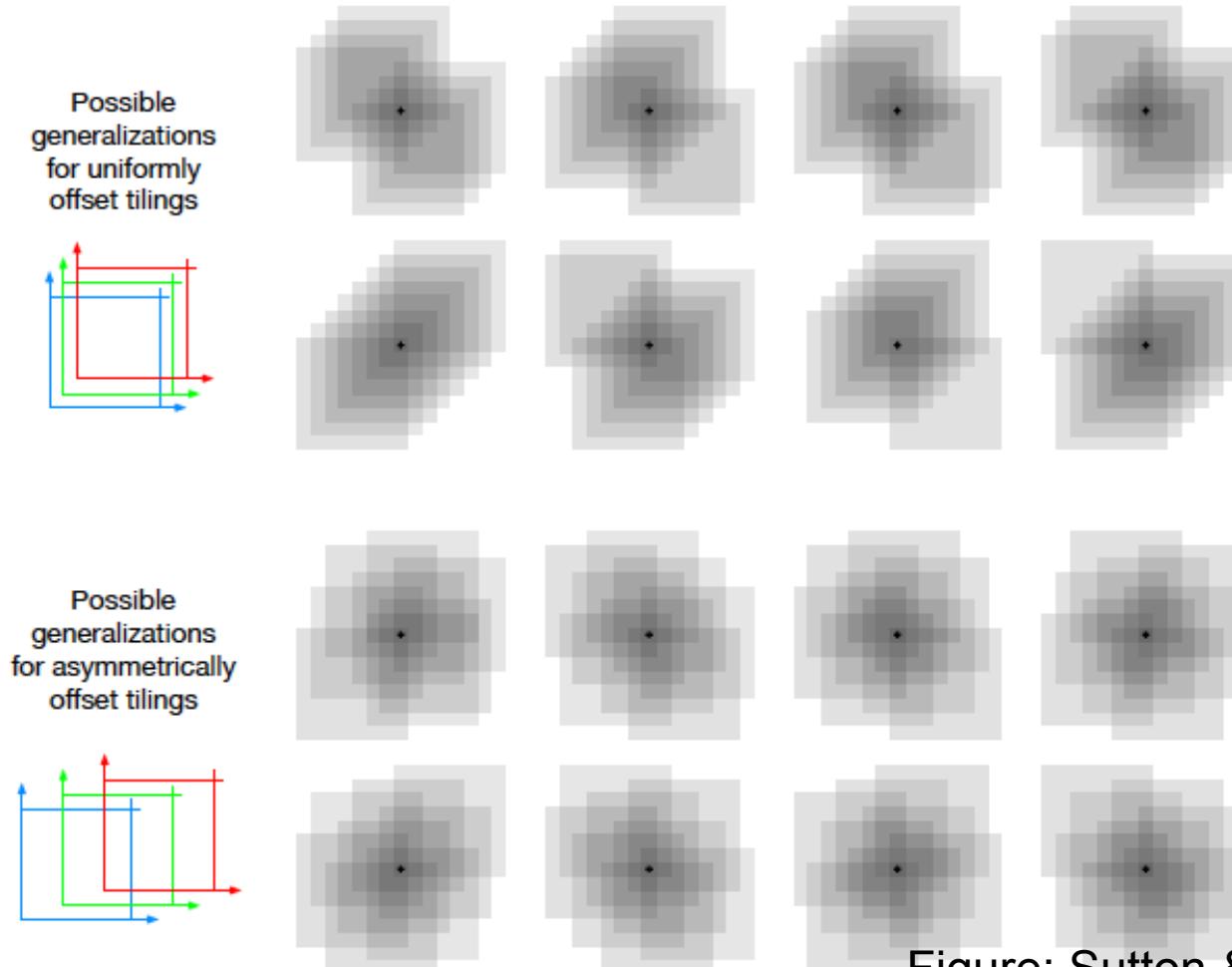
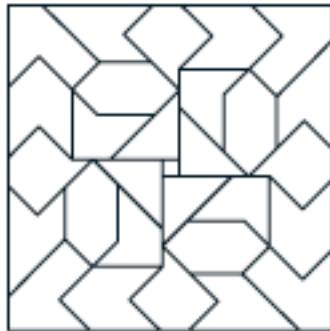


Figure: Sutton & Barto. RL:AI

Extra: Aggregation and tiling

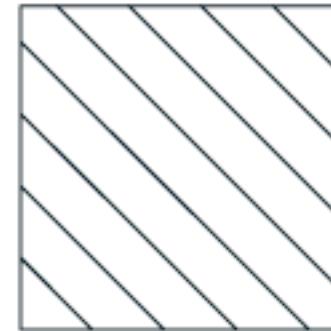
Special aggregations



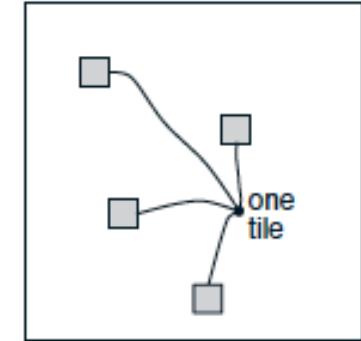
Irregular



Log stripes



Diagonal stripes



Hashing

Figure: Sutton & Barto. RL:AI

Generic approaches: Fourier basis

In multi-d, consider interaction between original dimensions
 n : max frequency in any direction; k : # original dimensions

$$x_i(s) = \cos(\pi s^\top \mathbf{c}^i)$$

$$\mathbf{c}^0 = [0, 0, \dots]^T$$

$$\mathbf{c}^1 = [1, 0, \dots]^T$$

$$\mathbf{c}^2 = [0, 1, \dots]^T$$

⋮

$$\mathbf{c}^{(n+1)^k} = [n, n, \dots]^T$$

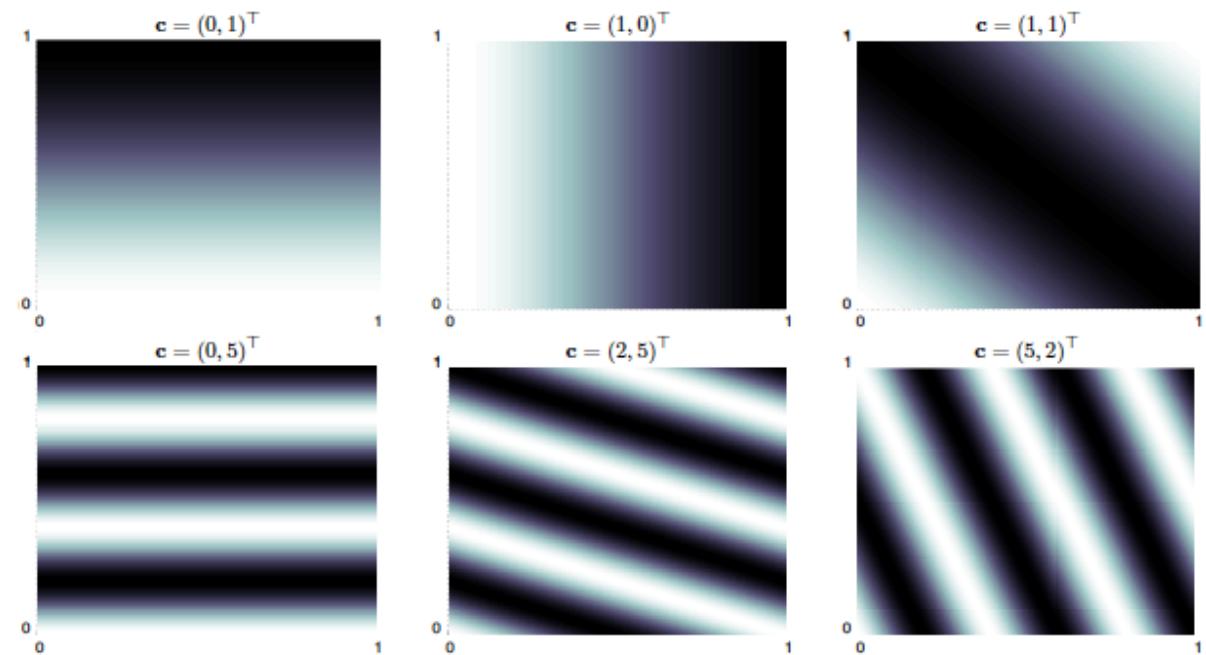


Figure: Sutton & Barto. RL:AI

Extra: Non-parametric approximation

Another way to approximate value functions is to directly use the training data rather than a parametric function

- Nearest neighbour
- Weighted average of n-nearest neighbour
- Locally weighted regression

More flexible, more precise where there is a lot of data

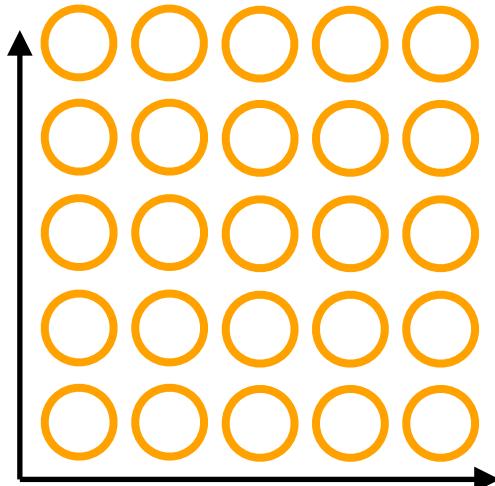
However, finding the nearest neighbours at every step is usually more expensive than evaluating a parametric function

Extra: Non-parametric approximation

Radial basis functions

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$

predefined features



Stationary kernel

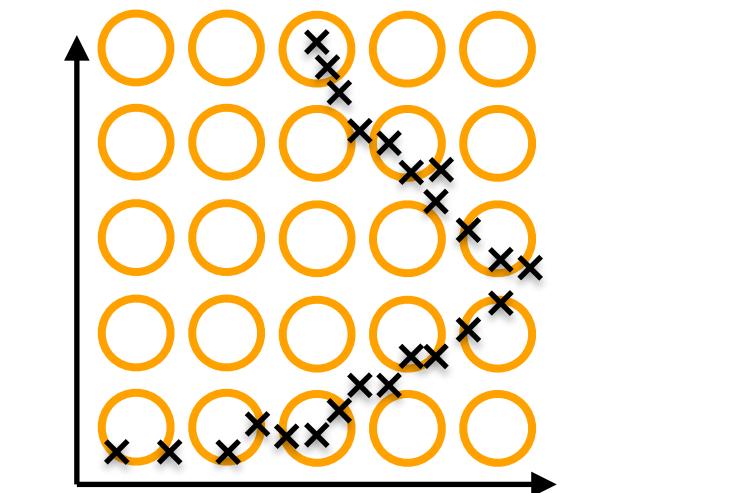
$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{s}_i\|^2}{2\sigma_i^2}\right)$$

data points

Extra: Non-parametric approximation

Radial basis functions

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$



Stationary kernel

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{s}_i\|^2}{2\sigma_i^2}\right)$$

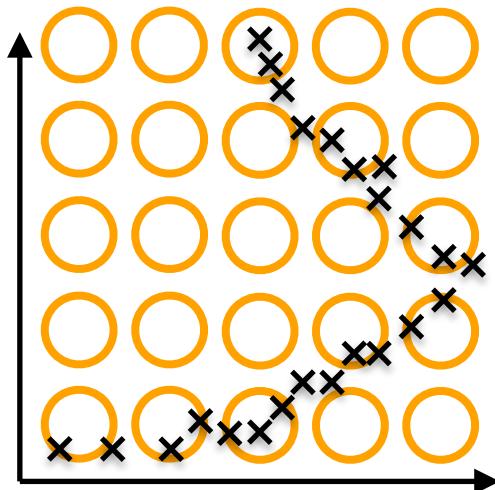
data points

Extra: Non-parametric approximation

Radial basis functions

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$

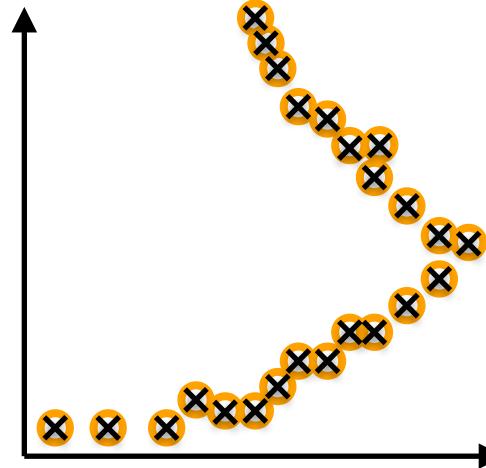
predefined features



Stationary kernel

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{s}_i\|^2}{2\sigma_i^2}\right)$$

data points

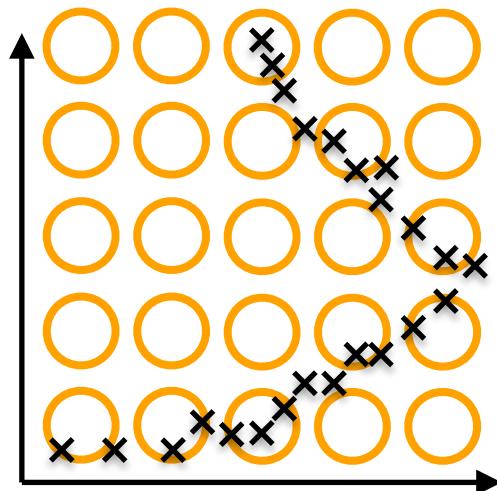


Extra: Non-parametric approximation

Radial basis functions

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$

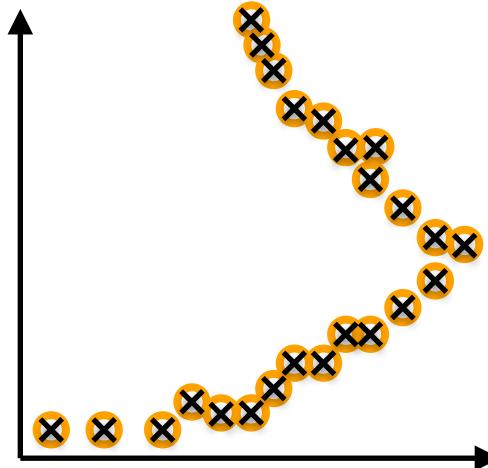
predefined features



Stationary kernel

$$\hat{v}(\mathbf{s}, \mathbf{w}) = \sum_i w_i \exp\left(-\frac{\|\mathbf{s} - \mathbf{s}_i\|^2}{2\sigma_i^2}\right)$$

data points



The kernel based approach has many of the nice properties of the linear function approximation case, while being highly flexible (more info in ML1)