# Recap

# Today's plan

First: Recap & Exam info

Time left: Q&A.

# Exam: General tips

Focus on understanding the methods and the relationship between them rather than on remembering e.g. update equations

Especially important: know the advantages, disadvantages and limitations of each methods, and the situations where a certain method should be preferred.

The recap will try to sketch the coherence of all lecture topics, nevertheless, we cannot cover all 14 lectures in 90 minutes, and topics outside of the recap can be on the exam, too.

# Exam: Cheat sheet

Many algorithms have variants:
Q- and V- version
importance weights
etc

Cheatsheet has most important ones only... Will be on Canvas after lecture to familiarise
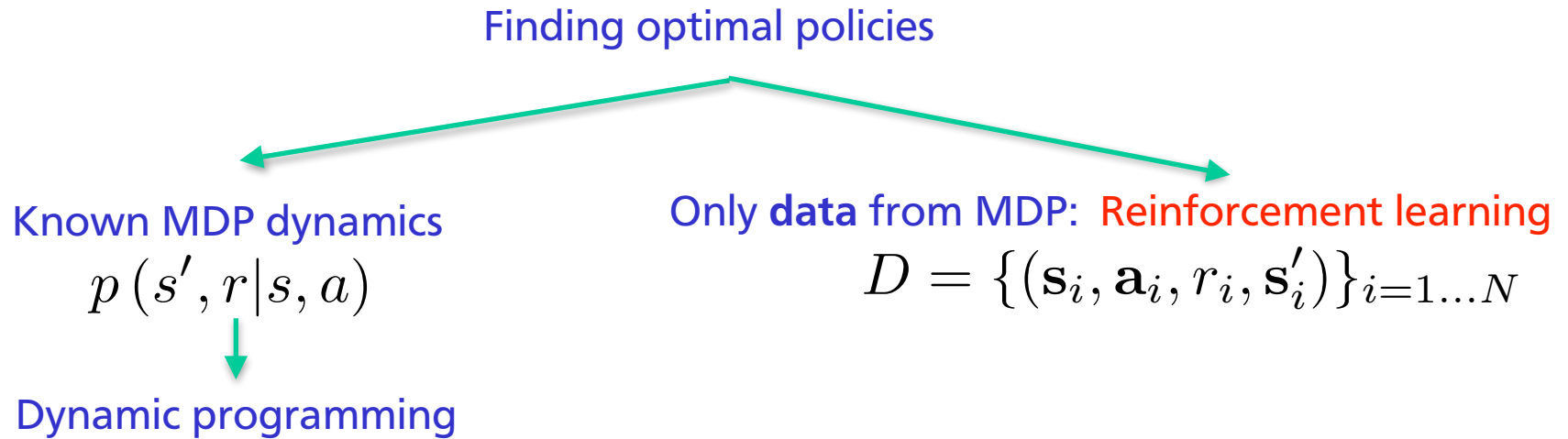
**Update equations cheat sheet**

- DP value iteration: $v_{k+1}(s) = \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')]$
- DP policy evaluation: $v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')]$
- Monte Carlo: $V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$
- TD(0): $V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$
- SARSA: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$
- Expected SARSA: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)]$
- Q-learning: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$
- $n$-step TD: $V_{t+n}(S_t) \doteq V_{t+n-1}(S_t) + \alpha[G_{t:t+n} - V_{t+n-1}(S_t)]$
- Tree backups: $Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha[G_{t:t+n}^{\text{tree}} - Q_{t+n-1}(S_t, A_t)]$, where $G_{t:t+n}^{\text{tree}} \doteq R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+n-1}(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1}) G_{t+1:t+n}^{\text{tree}}$
- Gradient Monte Carlo: $\mathbf{w} \leftarrow \mathbf{w} + \alpha[G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$
- Semi-gradient TD: $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})$
- LSTD: $\mathbf{w}_t \doteq \widehat{\mathbf{A}}_t^{-1} \widehat{\mathbf{b}}_t$, where $\widehat{\mathbf{A}}_t \doteq \sum_{k=0}^{t-1} \mathbf{x}_k (\mathbf{x}_k - \gamma \mathbf{x}_{k+1})^\top + \varepsilon \mathbf{I}$    and    $\widehat{\mathbf{b}}_t \doteq \sum_{k=0}^{t-1} R_{k+1} \mathbf{x}_k$
- GTD2: $\mathbf{v}_{t+1} \doteq \mathbf{v}_t + \beta \rho_t (\delta_t - \mathbf{v}_t^\top \mathbf{x}_t) \mathbf{x}_t, w_{t+1} = \mathbf{w}_t + \alpha \rho_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1}) \mathbf{x}_t^\top \mathbf{v}_t$
  (note: in the GTD2 equation, v is a parameter vector, not the value function)

Note: the following methods are given assuming the discount factor $\gamma = 1$.

- Finite difference gradients: $\theta_{k+1} = \theta_k + \alpha \frac{J(\theta_k - \epsilon) - J(\theta_k + \epsilon)}{2\epsilon}$
- REINFORCE: $\theta_{k+1} = \theta_k + \alpha G(\tau) \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t)$
- G(PO)MDP: $\theta_{k+1} = \theta_k + \alpha \sum_{t=0}^T r_t \sum_{t'=0}^t \nabla_\theta \log \pi_\theta(a_t|s_t)$
- PGT Actor-Critic: $\theta_{t+1} = \theta_t + \alpha \hat{q}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t|s_t)$
- Deterministic policy gradient (DPG): $\theta_{t+1} = \theta_t + \alpha \nabla_\theta \pi_\theta(a_t|s_t) \nabla_a q(s_t, a_t)|_{a=\pi_\theta(s)}$
- Natural policy gradient: $\theta_{t+1} = \theta_t + \alpha F^{-1}(\theta) \nabla_\theta J(\theta)$
  where $\nabla_\theta J(\theta)$ is an estimate of the 'vanilla' policy gradient.

# Big picture

Finding optimal policies

Known MDP dynamics
$$p\left(s', r | s, a\right)$$

Dynamic programming

Only **data** from MDP:  Reinforcement learning
$$D = \left\{\left(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\right)\right\}_{i=1...N}$$

Planning Method.

Value Iteration

$$V_t^{\pi}(s) = \max_a \mathbb{E}_{s'}\left[r(s,a) + \gamma V_{t+1}^{\pi}(s')\right]$$



Policy Iteration.

1) Policy Evaluation

until convergen $V_{\pi}$

2) $\pi_{new}(s) = \text{argmax} \, \mathbb{E}_{s'}\, r(s,a) + \gamma V_{\pi}(s')$

# Big picture: How to learn policies

$$D = \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}_{i=1...N}$$

Learn value function
$$V(\mathbf{s}), Q(\mathbf{s}, \mathbf{a})$$

Learn model
$$p(\mathbf{s}'|\mathbf{s}, \mathbf{a}) \ r(\mathbf{s}, \mathbf{a})$$

Policy
$$\pi(\mathbf{a}|\mathbf{s})$$

Optimize policy
$$\pi(\mathbf{a}|\mathbf{s})$$

Optimize policy
$$\pi(\mathbf{a}|\mathbf{s})$$

Learn value or policy
$$V(\mathbf{s}), Q(\mathbf{s}, \mathbf{a}) \ \pi(\mathbf{a}|\mathbf{s})$$

Policy
$$\pi(\mathbf{a}|\mathbf{s})$$

**Critic only**

**Actor-critic**

**Actor only**

**Model-free RL**

**Model-based RL**

*Transition model*

Thanks to Jan Peters

# Value-based methods: MC

(Critic - only)
Value fc methods → TD

L 2,3

MC

L 5

tabular MC

gradient MC

$$V_\pi(S_t) \leftarrow V_\pi(S_t) + \alpha \left( G_t - V_\pi(S_t) \right)$$

'target'

mean-squared value error

$$\sum_S \mu(s) \left( V(s) - \hat{V}_w(s) \right)^2$$
$$\pi$$

$$w \leftarrow w + \alpha \left( G_t - \hat{V}_w(s) \right) \cdot \nabla_w \hat{V}_w(s)$$

+ imp weights.

Herke van Hoof |    8

Reinforcement Learning

# Value based methods: TD learning

on-policy

TD-Learning

off-policy

$$V(S_t) \leftarrow V(S_t) + \alpha(\text{target} - V(S_t))$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{target} - Q(S_t, A_t))$$

TD(0),    SARSA

expected SARSA

Q-learning

$$\text{target} = R_{t+1} + \gamma V(S_{t+1}); \; R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$$

target   $R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) \cdot Q(S_{t+1}, a)$

$$R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$$

n-step TD   /   n-step SARSA

Semi-gradient TD(0)

$$W \leftarrow W + \alpha(R + \gamma \hat{V}_w(S') - \hat{V}_w(S)) \cdot \nabla_w \hat{V}_w(S)$$

GTD2
gradient of the
ms projected BE

double Q
learning

DQN

- semi-gradient Q-learning
- experience replay,
  target networks.

$$R_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$
$$+ \gamma^n Q(S_{t+n}, A_{t+n})$$

(+imp. wts)

# Off-policy learning

| On-policy | Off-policy |
|-----------|-----------|
| Simpler | More complex |
| Specific case | More general |
| Often converges faster | Often large variance or slow convergence |
| Only data gathered with current policy | Can reuse data, use data from other source |
| Generally needs non-greedy policy | Allows greedy target policy |

# Evaluation methods

distribution model
or real/simulated samples

distribution model

distribution model
or real/simulated samples

distribution model

Figure from
Sutton and Barto RL:AI

Reinforcement learning

# Evaluation methods

distribution model
or real/simulated samples

distribution model

Temporal-difference learning

width of update

Dynamic programming

depth (length) of update

N-Step

gradient
Monte-Carlo

distribution model
or real/simulated samples

Monte Carlo

Exhaustive search

distribution model

Figure from
Sutton and Barto RL:AI

Herke van Hoof | 11

Reinforcement learning

# Evaluation methods

GTD2, semi-gradient TD(0) & LSTD

distribution model or real/simulated samples

Temporal-difference learning

width of update

Dynamic programming

distribution model

depth (length) of update

N-Step

gradient Monte-Carlo

distribution model or real/simulated samples

Monte Carlo

Exhaustive search

distribution model

Figure from Sutton and Barto RL:AI

Herke van Hoof | 11

# Control methods

On-policy
(narrow)

Sarsa

N-step
TD

Monte-
Carlo
control

Depth

Off-policy
(wide)

Q-learning &
Expected SARSA

$(q_*)$    $s, a$

$s, a$

$r$

$p$   $r$

$s'$

$s'$

DP
(wider)

max

$\pi$

$a'$

$a'$

Policy evaluation &
Value iteration

# Types of function approximation

For any of the methods (gradient MC / semi-gradient TD/ LSTD / GTD2), choice of function approximation

| linear | non-linear |
|---|---|
| tabular | e.g. neural network |
| aggregate | |
| tiling | |
| radial basis function | |
| polynomial basis function | |
| fourier basis function | |

# Convergence with function approximation

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

**Local or global convergence?**

* with appropriate step-size schedule          ** if features independent, single solution

Reinforcement Learning

# Convergence with function approximation

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

**Convergence to minimum of which error?**

\* 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 * | global | global | non-linear: local | MC:VE | global | non-linear: local |
| Semi-gradient TD * | | | | TD: PBE | 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

# Semi-gradients?

Semi-gradient $TD(0)$

$$w \leftarrow w + \gamma \left( R + \hat{V}_w(s') - \hat{V}_w(s) \right) \cdot \nabla_w \hat{V}_w(s)$$

not gradient of any thing.

looks like gradient
mS TD error.
↓
Converge to minimum TD error.
↓
Weird result, depency of value
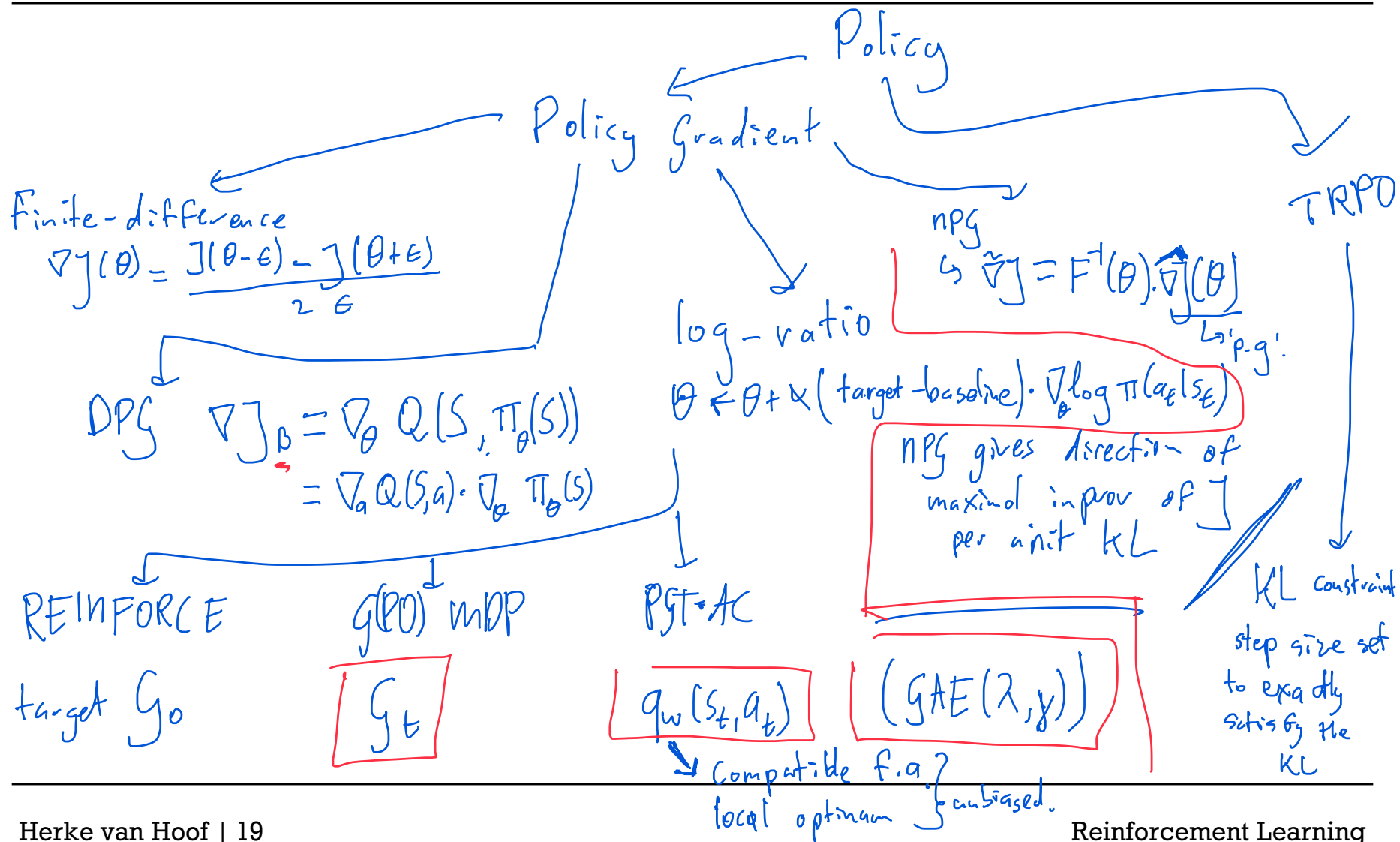     on the past
↓
we don't do this, typically.

→ TD fix point ─────┐
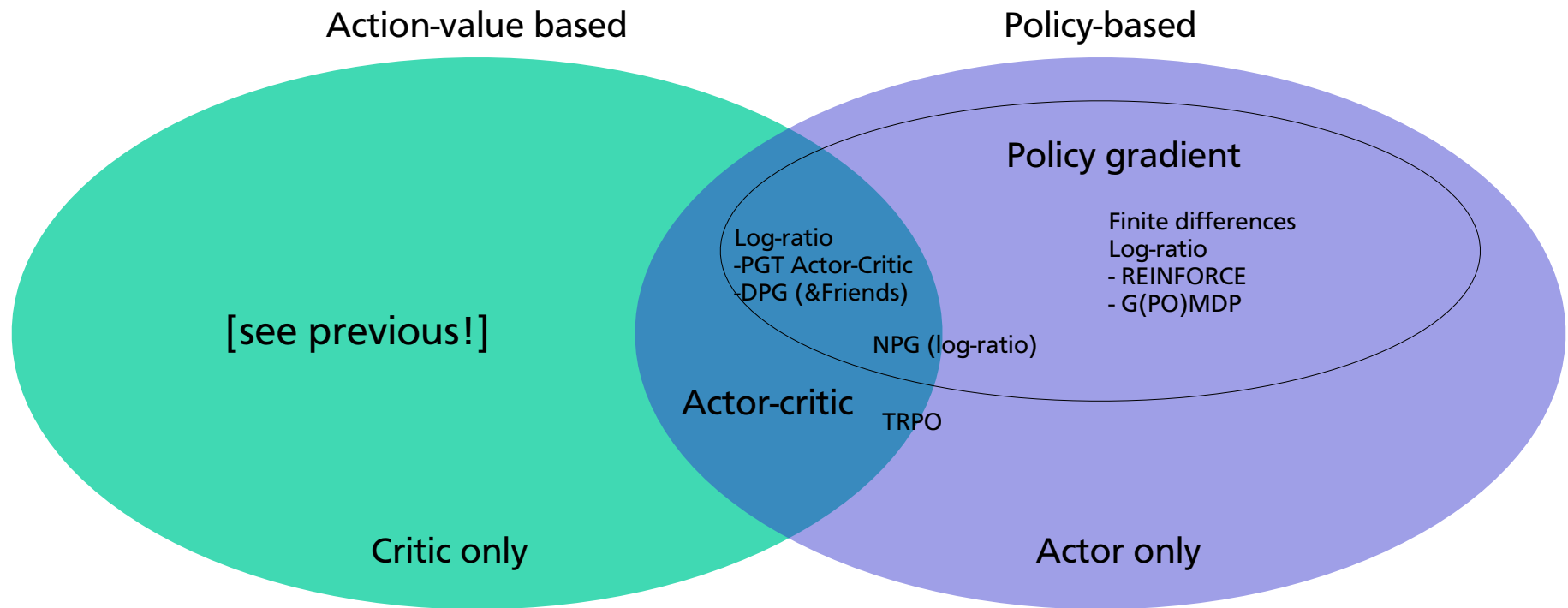        ↓            found by
      MSPBE          LSTD.
        ↓ gradient
      GTD2

Policy

Policy Gradient

Finite-difference

$$\nabla J(\theta) = \frac{J(\theta-\epsilon) - J(\theta+\epsilon)}{2\epsilon}$$

DPG

$$\nabla J_\beta = \nabla_\theta Q(S, \pi_\theta(S))$$
$$= \nabla_a Q(S,a) \cdot \nabla_\theta \pi_\theta(S)$$

REINFORCE

target $G_0$

GPO) MDP

$$\boxed{G_t}$$

log-ratio

$$\theta \leftarrow \theta + \alpha (\text{target} - \text{baseline}) \cdot \nabla_\theta \log \pi(a_t | s_t)$$

PGT·AC

$$\boxed{q_w(s_t, a_t)}$$

compatible f.a?
local optimum {unbiased.

nPG

$$\leftrightarrow \tilde{\nabla} J = F^{-1}(\theta) \cdot \widehat{\nabla J}(\theta)$$
$\hookrightarrow$ p.g.

nPG gives direction of
maximal improv of J
per unit KL

$$\boxed{(GAE(\lambda, \gamma))}$$

TRPO

KL constraint
step size set
to exactly
satisfy the
KL

# Policies and action-values

Action-value based

Policy-based

Policy gradient

[see previous!]

Log-ratio
-PGT Actor-Critic
-DPG (&Friends)

Finite differences
Log-ratio
- REINFORCE
- G(PO)MDP

NPG (log-ratio)

Actor-critic

TRPO

Critic only

Actor only

Reinforcement Learning

# Model-based learning

*Transition-model.*

- Types of models (generative, trajectory, distributional)

- Dyna-Q

  • Prioritized sweeping, on-policy & uniform sampling

- Planning a full policy or from current state only (See lecture 11)

- Backpropagation through the model

# State update functions in POMDPs

Why do we need update functions for internal states?

What are the properties of int. states (compactness, markovian)

## Exact methods

- **Full history**
  Not compact...

- **Belief state**
  Easy to interpret *(Comp. heavy)*
  Requires known model

- **Predictive state**
  Model learnable from data
  Most compact

## Approximate methods

- **Recent observation(s)**
  Easy
  Lose long-term dependencies

- **End-to-end learning**
  Quite general
  RNN learning can be tricky,
  requires much data...

# Other topics that are important

<mark>Maximization bias</mark> (lecture 4)

Exploration vs exploitation (throughout, also lecture 1)

Pure exploration & best arm identification (lecture 13)

Anything else I missed?

# Other tips

==Very important== - know ==when/why to use each method/strategy== ==(advantages, disadvantages&limitations)==

**Good luck with your preparation!**

# Intentionally left blank