# State Abstractions
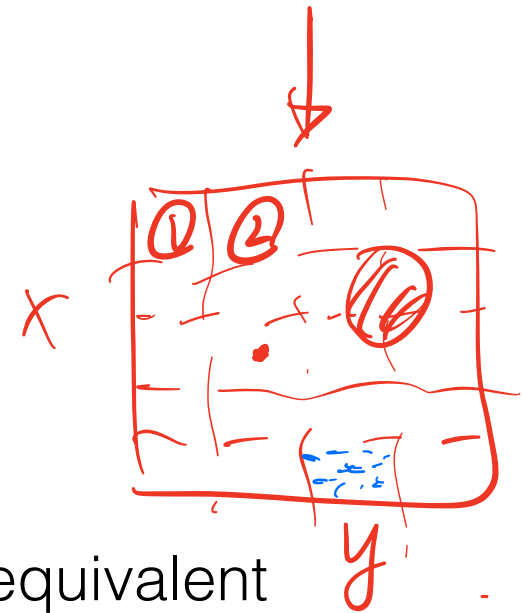
# Notations and Setup

- MDP $M = (S, A, P, R, \gamma)$

- Abstraction $\phi : S \to S_\phi$

- Surjection — aggregate states and treat as equivalent

- Are the aggregated states really equivalent?

- Do they have the same…

  - optimal action?

  - Q* values?

  - dynamics and rewards?

$$(x, y, z, u)$$
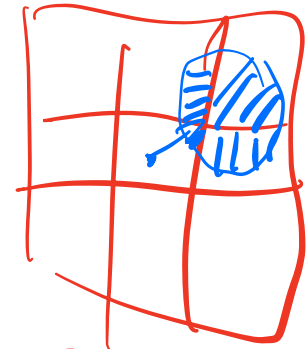
$$S$$
$$\|$$
$$(x, y, z, w) \overset{\phi}{\longmapsto} (x, y)$$

# Outline of the lecture

1. Define various notions/criteria of abstractions

2. Study their relationships

3. Analyze how to use them (e.g., building an abstract model) in planning and learning

   - Abstract model will also appear in 1 & 2

# Abstraction hierarchy

An abstraction $\phi$ is … if … $\forall s^{(1)}, s^{(2)}$ where $\phi(s^{(1)}) = \phi(s^{(2)})$

- $\pi^*$-irrelevant: $\exists \pi_M^*$ s.t. $\pi_M^*(s^{(1)}) = \pi_M^*(s^{(2)})$

- $Q^*$-irrelevant: $\forall a$, $Q_M^*(s^{(1)}, a) = Q_M^*(s^{(2)}, a)$

- Model-irrelevant: $\forall a \in A$, $\qquad\qquad R(s^{(1)}, a) = R(s^{(2)}, a)$
  (bisimulation)
  $\qquad\qquad \forall a \in A, x' \in S_\phi, \quad P(x' \mid s^{(1)}, a) = P(x' \mid s^{(2)}, a)$

$$\sum_{s' \in \phi^{-1}(x')} P(s' \mid s^{(1)}, a)$$

$\forall s' \in S.$

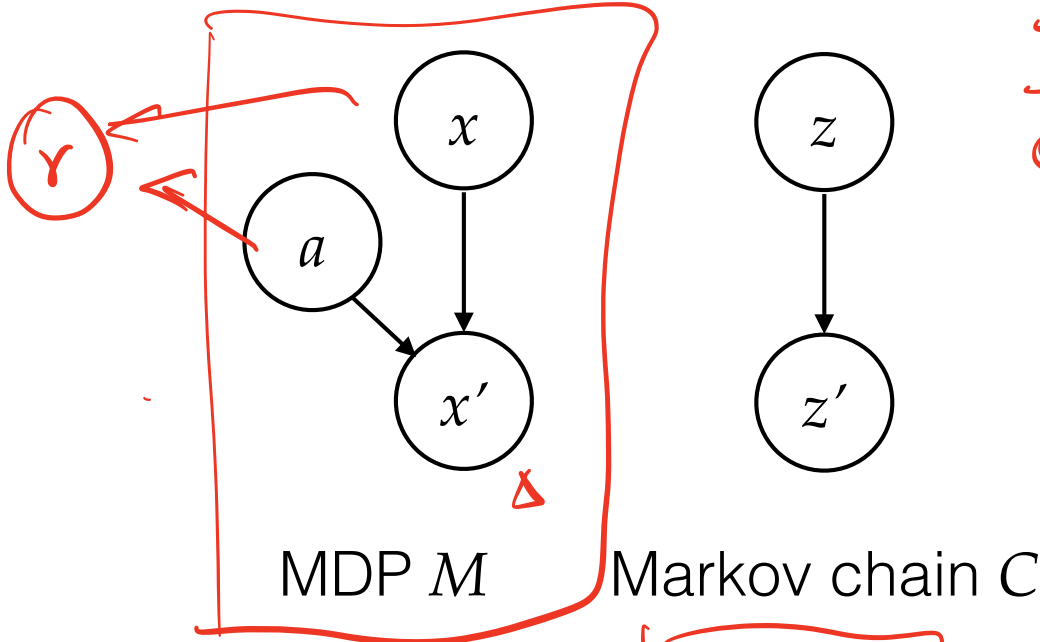$P(s' \mid s^{(1)}, a)$

$= P(s' \mid s^{(2)}, a)$

**Theorem:** Model-irrelevance $\Rightarrow Q^*$-irrelevance $\Rightarrow \pi^*$-irrelevance

Why not $P(s' \mid s^{(1)}, a) = P(s' \mid s^{(2)}, a)$ ?

violate.

$S = (x, z)$.

$\phi : (x, z) \mapsto \boxed{x}$.



MDP $M$  Markov chain $C$

$(x, z^{(1)})$ and $(x, z^{(2)})$ cannot be aggregated under the $s'$-based condition

$P((x', z') \mid (x, z), a) = P_M(x' \mid x, a) \cdot P_C(z' \mid z)$

$(x, z^{(1)})$

$(x, z^{(2)})$.

integrated out by bisimulation

# Abstraction induces an **equivalence relation**

- Reflexivity, symmetry, transitivity

- Equivalence notion is a canonical representation of abstraction
  (i.e., what symbol you associate with each abstract state doesn't matter; what
  matters is which states are aggregated together)

- Partition the state space into **equivalence classes**

- Coarsest bisimulation is unique (see proof in notes)

  - sketch: if $\phi_1$ and $\phi_2$ are both bisimulations, their *common
    coarsening* is also a bisimulation (two states are aggregated if
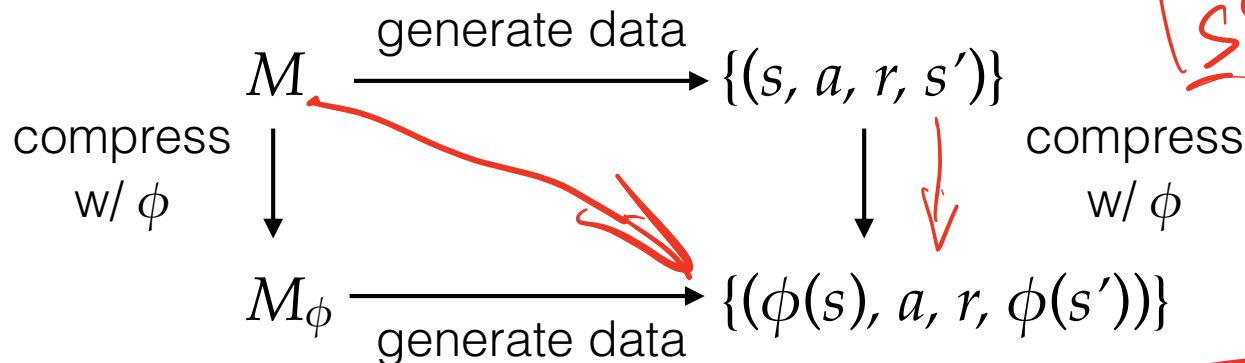    they are aggregated under *either* $\phi_1$ or $\phi_2$)

# The abstract MDP implied by bisimulation

$\phi$ is bisimulation: $R(s^{(1)}, a) = R(s^{(2)}, a)$, $P(x' \mid s^{(1)}, a) = P(x' \mid s^{(2)}, a)$

$\in S_\phi$

- MDP $M_\phi = (S_\phi, A, P_\phi, R_\phi, \gamma)$

- For any $x \in S_\phi$, $a \in A$, $x' \in S_\phi$

  - $R_\phi(x, a) = R(s, a)$ for any $s \in \phi^{-1}(x)$

  - $P_\phi(x' \mid x, a) = P(x' \mid s, a)$ for any $s \in \phi^{-1}(x)$

$(s, a, r, s')$

$\downarrow$

$(\phi(s), a, r, \phi(s'))$

- No way to distinguish between the two routes:

$s^{(1)}, a \quad \# \, n.$
$s^{(2)}, a \sim n.$

$x, a \sim 2n.$

$$M \xrightarrow{\text{generate data}} \{(s, a, r, s')\}$$

compress w/ $\phi$ $\downarrow$ $\qquad$ compress w/ $\phi$ $\downarrow$

$$M_\phi \xrightarrow[\text{generate data}]{} \{(\phi(s), a, r, \phi(s'))\}$$

7

# Implications of bisimulation

- $Q^*$ is preserved

- $Q_M^\pi$ is preserved *for any $\pi$ lifted from an abstract policy*

  - the policy must take the same action (distribution) across aggregated states

$$\bar{\pi}: S_\phi \to (A)$$

$$s \longmapsto \pi(\phi(s))$$

$$\Rightarrow R(s^{(1)}, a) = R(s^{(2)}, a)$$

$$\to P(x' \mid s^{(1)}, a) = P(x' \mid s^{(2)}, a)$$

8

# Extension to handle action aggregation

$$(S, A, P, R, \gamma)$$



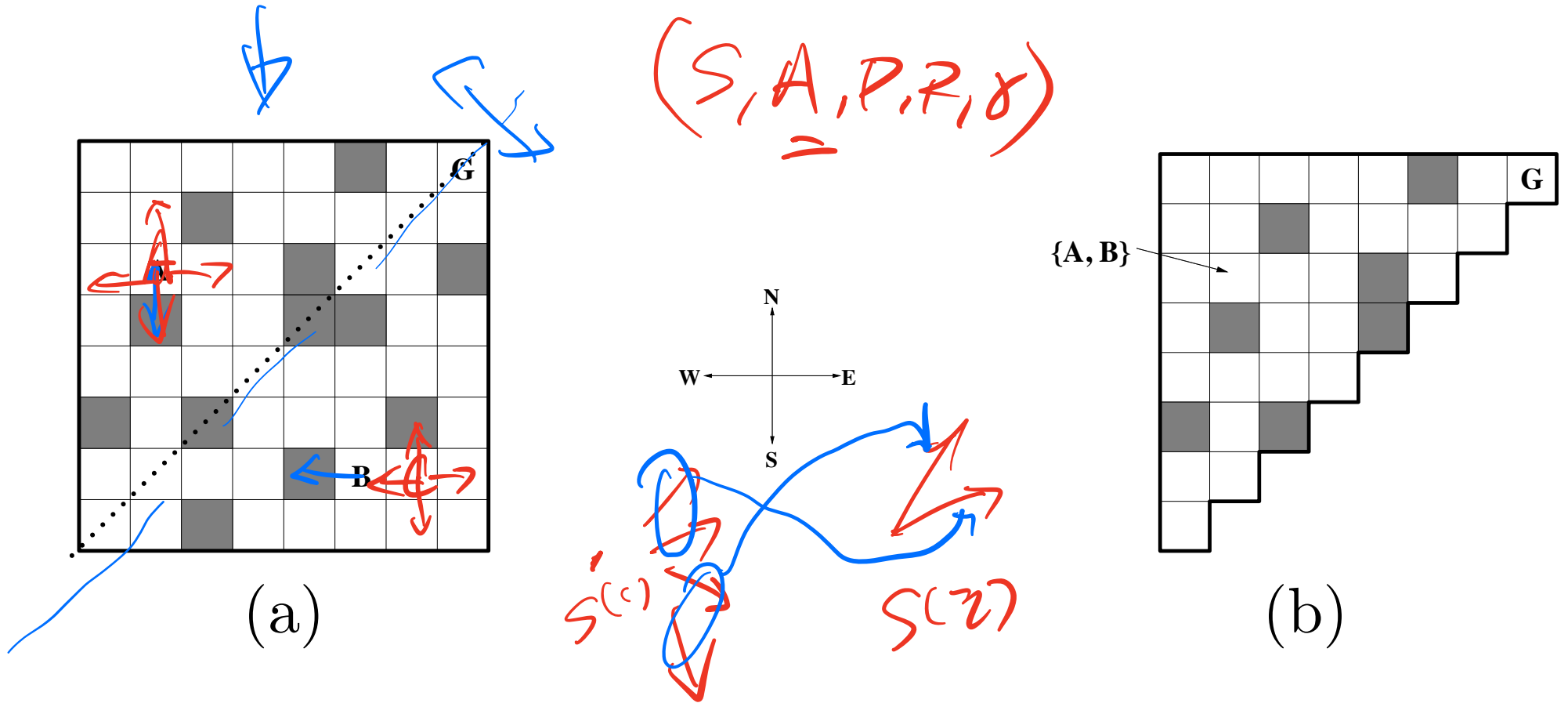(a)

N
W — E
S

$S^{(1)}$  $S^{(2)}$

{A, B}

G

(b)

Figure from: Ravindran & Barto. Approximate Homomorphisms: A framework for non-exact minimization in Markov Decision Processes. 2004.

$[\pi]_M$

**Definition 3** (Approximate abstractions). Given MDP $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$ and state abstraction $\phi$ that operates on $\mathcal{S}$, define the following types of abstractions:

1. $\phi$ is an $\epsilon_{\pi^*}$-approximate $\pi^*$-irrelevant abstraction, if there exists an abstract policy $\pi : \mathcal{S}_\phi \to \mathcal{A}$, such that $\|V_M^* - V_M^{[\pi]_M}\|_\infty \leq \epsilon_{\pi^*}$.

2. $\phi$ is an $\epsilon_{Q^*}$-approximate $Q^*$-irrelevant abstraction if there exists an abstract $Q$-value function $f : \mathcal{S}_\phi \times \mathcal{A} \to \mathbb{R}$, such that $\|[f]_M - Q_M^*\|_\infty \leq \epsilon_{Q^*}$.

3. $\phi$ is an $(\epsilon_R, \epsilon_P)$-approximate model-irrelevant abstraction if for any $s^{(1)}$ and $s^{(2)}$ where $\phi(s^{(1)}) = \phi(s^{(2)})$, $\forall a \in \mathcal{A}$,

$$P(x' \mid s^{(i)}, a)$$

$$|R(s^{(1)}, a) - R(s^{(2)}, a)| \leq \epsilon_R, \quad \left\|\Phi P(s^{(1)}, a) - \Phi P(s^{(2)}, a)\right\|_1 \leq \epsilon_P. \tag{3}$$

Useful notation: $\Phi$ is a $|\mathcal{S}_\phi| \times |\mathcal{S}|$ matrix, with

$$\Phi(x, s) = \mathbb{I}[\phi(s) = x]$$

- lifting a state-value function: $[V_{M_\phi}^\star]_M = \Phi^\top V_{M_\phi}^\star$

- collapsing the transition distribution: $\Phi P(s, a)$

10

$$\Phi = \begin{bmatrix} | & | & | & & & & \\ & & | & | & | & | & \\ & & & & \ddots & & \\ & & & & & & | & | \end{bmatrix}$$

(1) $\quad p \in \Delta(S) \subseteq \mathbb{R}^{|S|}$

$$\Phi \times p = \begin{bmatrix} | & | & | & & & \phi \\ & & | & | & & \\ & & & \ddots & & \\ & & & & |\phi| \end{bmatrix} \begin{bmatrix} \vdots \\ p \\ ( \end{bmatrix}$$

$$\Phi\, P(\cdot \mid s^{(1)}, a) = \Phi\, P(\cdot \mid s^{(2)}, a)$$

(2) $\quad \boxed{f} : S_\phi \to \mathbb{R}. \quad \Phi^T f.$

$$\begin{bmatrix} \rule{2cm}{0.4pt} \\ \rule{2cm}{0.4pt} \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} 0 \\ f \\ \vdots \end{bmatrix} = [f]_M . \begin{bmatrix} \vdots \; 0 \leftarrow \\ 0 \leftarrow \\ 0 \leftarrow \end{bmatrix}$$

**Theorem 2.** *(1) If $\phi$ is an $(\epsilon_R, \epsilon_P)$-approximate model-irrelevant abstraction, then $\phi$ is also an approximate $Q^\star$-irrelevant abstraction with approximation error $\epsilon_{Q^\star} = \frac{\epsilon_R}{1-\gamma} + \frac{\gamma\epsilon_P R_{\max}}{2(1-\gamma)^2}$.*
*(2) If $\phi$ is an $\epsilon_{Q^\star}$-approximate $Q^\star$-irrelevant abstraction, then $\phi$ is also an approximate $\pi^\star$-irrelevant abstraction with approximation error $\epsilon_{\pi^\star} = 2\epsilon_{Q^\star}/(1-\gamma)$.*

$$\left(\|V^\star - V^{\pi_f}\|_\infty \le \frac{2\|f - Q^\star\|}{1-\gamma}\right)$$

- (2) follows directly from a known result; can you see?

- Construct the $f$ in the definition of approx. $Q^\star$-irrelevance:

$\phi$ is an $\epsilon_{Q^\star}$-approximate $Q^\star$-irrelevant abstraction if there exists an abstract $Q$-value function $f : S_\phi \times A \to \mathbb{R}$, such that $\|[f]_M - Q^\star_M\|_\infty \le \epsilon_{Q^\star}$.

- Define $M_\phi = (S_\phi, A, P_\phi, R_\phi, \gamma)$ w/ any weighting distributions $\{p_x : x \in S_\phi\}$, where each $p_x$ is supported on $\phi^{-1}(x)$

$$R_\phi(x, a) = \Sigma_{s \in \phi^{-1}(x)}\, p_x(s)\, R(s, a), \quad P_\phi(x, a) = \Sigma_{s \in \phi^{-1}(x)}\, p_x(s)\, \Phi\, P(s, a).$$

- $\left|R_\phi(\phi(s), a) - R(s, a)\right| \le \varepsilon_R, \quad \left\|P_\phi(\phi(s), a) - \Phi\, P(s, a)\right\| \le \varepsilon_P.$

- Set $f := Q^\star_{M_\phi}$, bound $\|[f]_M - Q^\star_M\|_\infty$

$$\left\| \left[ Q^*_{M\phi} \right]_M - Q^*_M \right\|_\infty$$

$$= \frac{1}{1-\gamma} \left\| \left[ Q^*_{M\phi} \right]_M - \mathcal{T}_M \left[ Q^*_{M\phi} \right]_M \right\|_\infty$$

$$\| g - Q^* \|_\infty$$
$$= \| g - \mathcal{T}g + \mathcal{T}g - \mathcal{T}Q^* \|_\infty$$
$$\leq \| g - \mathcal{T}g \|_\infty + \gamma \| g - Q^* \|_\infty$$

$$= \frac{1}{1-\gamma} \left\| \left[ \mathcal{T}_{M\phi} Q^*_{M\phi} \right]_M - \mathcal{T}_M \left[ Q^*_{M\phi} \right]_M \right\|_\infty.$$

$$\forall (s,a).$$

$$\left| \left( \mathcal{T}_{M\phi} Q^*_{M\phi} \right) (\phi(s), a) - \mathcal{T}_M \left[ Q^*_{M\phi} \right]_M (s,a) \right|$$

$$= \left| R_\phi(\phi(s), a) + \gamma \langle P_\phi(\phi(s), a), V^*_{M\phi} \rangle - R(s,a) + \gamma \langle P(s,a), [V^*_{M\phi}]_M \rangle \right| \quad R(s,a) + \gamma \mathbb{E}_{s \sim P(\cdot | s,a)} \left[ \max_{a'} Q(s',a') \right].$$

$$\leq \varepsilon_R + \gamma * \qquad s' \to \max_{a'} Q(s',a')$$

$$\left( \langle P_\phi(\phi(s), a), V^*_{M\phi} \rangle - \langle P(s,a), \Phi^\top V^*_{M\phi} \rangle \right).$$

$$\varepsilon_p \cdot \frac{V_{max}}{2}. \qquad \langle \Phi P(s,a), V^*_{M\phi} \rangle \quad \langle u, Av \rangle := (u^\top A)v$$

$$= \langle A^\top u, v \rangle.$$

$$\Phi \ P(s^{(1)}, a) \ = \ \Phi \ P(s^{(2)}, a).$$



$$\mathbb{E}_p[f] \ \neq \ \mathbb{E}_q[f].$$

# Outline of the lecture

1. Define various notions/criteria of abstractions

2. Study their relationships

3. Analyze how to use them (e.g., building an abstract model) in planning and learning

- e.g., plan in $M_\phi$ to reduce computational cost

- If $\phi$ is not exact bisimulation, what's sub-optimality as a function of $(\varepsilon_R, \varepsilon_P)$ ? (Partially answered; will take a closer look)

- What if $\phi$ is only approximately Q*-irrelevant? Is the abstract model still useful? Can we still bound loss as a function of $\varepsilon_{Q^*}$?

- Learning setting?

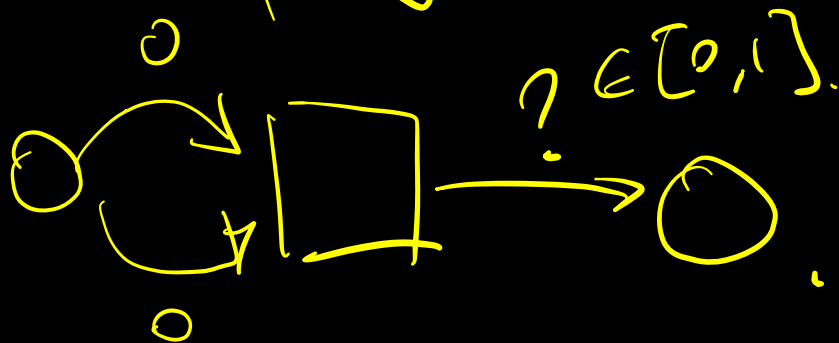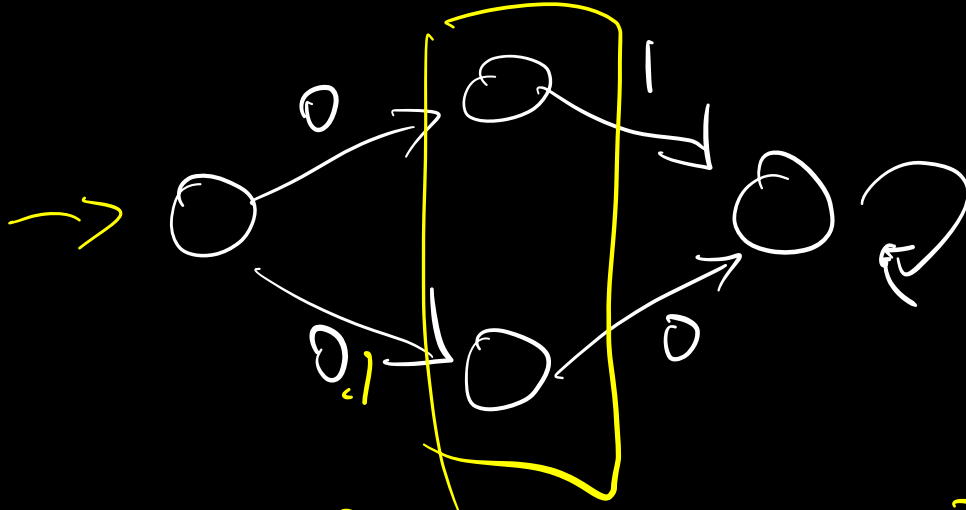# Loss of $\pi^\star_{M_\phi M}$: approx. bisimulation

- Recall: $M_\phi$ defined using any weighting distributions $\{p_x\}$ satisfies
  $|R_\phi(\phi(s), a) - R(s, a)| \le \varepsilon_R, \quad \|P_\phi(\phi(s), a) - \Phi\, P(s, a)\|_1 \le \varepsilon_P.$

- Apply earlier Theorem: $\left\| V_M^\star - V_M^{[\pi^\star_{M_\phi}]M} \right\|_\infty \le \frac{2\epsilon_R}{(1-\gamma)^2} + \frac{\gamma\epsilon_P R_{\max}}{(1-\gamma)^3}$

- Can improve: $\left\| V_M^\star - V_M^{[\pi^\star_{M_\phi}]M} \right\|_\infty \le \frac{2\epsilon_R}{1-\gamma} + \frac{\gamma\epsilon_P R_{\max}}{(1-\gamma)^2}$

- Idea: for any $\pi : S_\phi \to A$, $\quad \left\| [V_{M_\phi}^\pi]_M - V_M^{[\pi]M} \right\|_\infty \le \frac{\epsilon_R}{1-\gamma} + \frac{\gamma\epsilon_P R_{\max}}{2(1-\gamma)^2}$

- Finally,

$$V_M^\star(s) - V_M^{[\pi^\star_{M_\phi}]M}(s) = V_M^\star(s) - V_{M_\phi}^\star(\phi(s)) + V_{M_\phi}^\star(\phi(s)) - V_M^{[\pi^\star_{M_\phi}]M}(s)$$

$$\le \left\| Q_M^\star - [Q_{M_\phi}^\star]_M \right\|_\infty + \left\| [V_{M_\phi}^{\pi^\star_{M_\phi}}]_M - V_M^{[\pi^\star_{M_\phi}]M} \right\|_\infty$$

- Lesson: w/ approx. bisimulation, take the $\max_\pi \|V_M^\pi - V_{\widehat{M}}^\pi\|_\infty$ route instead of the $\|Q_M^\star - Q_{\widehat{M}}^\star\|$ route to save dependence on horizon

$[\pi^*_{M\phi}]_M$ can be bad. even if

$\phi$ is $\pi^*$-irrelevant.

# Loss of $\pi^\star_{M_{\phi M}}$ : approx. Q*-irrelevance

- $M_\phi$ well defined, but transitions/rewards don't make sense

- Can still show: $\|[Q^\star_{M_\phi}]_M - Q^\star_M\|_\infty \leq 2\epsilon_{Q^\star}/(1-\gamma)$

- Exact case ($\epsilon_{Q^\star} = 0$): $\forall\, s^{(1)}, s^{(2)}$ where $\phi(s^{(1)}) = \phi(s^{(2)})$

$$R(s^{(1)}, a) + \gamma\langle P(s^{(1)}, a), V^\star_M\rangle = Q^\star(s^{(1)}, a) = Q^\star(s^{(2)}, a) = R(s^{(2)}, a) + \gamma\langle P(s^{(2)}, a), V^\star_M\rangle$$

"inverse" of lifting (can only be applied to piece-wise constant functions)

So: $\quad (\mathcal{T}_{M_\phi}[Q^\star_M]_\phi)(x, a) = R_\phi(x, a) + \gamma\langle P_\phi(x, a), [V^\star_M]_\phi\rangle$

$$= \sum_{s\in\phi^{-1}(x)} p_x(s)\,(R(s, a) + \gamma\langle\Phi P(s, a), [V^\star_M]_\phi))$$

$$= \sum_{s\in\phi^{-1}(x)} p_x(s)\,(R(s, a) + \gamma\langle P(s, a), V^\star_M))$$

$$= \sum_{s\in\phi^{-1}(x)} p_x(s)\,[Q^\star_M]_\phi(x, a) = [Q^\star_M]_\phi(x, a).$$

$$\underline{[Q^*_{M\phi}]_M = Q^*_M} \qquad \rightarrow \text{ fixed point of } T_M.$$

$$\Longleftrightarrow \qquad Q^*_M \overset{\in \mathbb{R}^{S \times A}}{\text{(compressed)}} =: [Q_{M^*}]_{\underline{\phi}}.$$

$$\text{is fixed point of } \underline{T_{M\phi}}.$$

$$\left(\left(\boxed{T_{M\phi}}\right) [Q^*_M]_\phi\right)(x,a) \overset{?}{=} \boxed{[Q^*_M]_\phi (x,a).}$$

$$= \underline{R_\phi(x,a)} + \gamma \langle \underline{P_\phi(x,a)}, [V^*_M]_\phi \rangle$$

$$= \sum_{s \in \phi^{-1}(x)} \underline{P_x(s)} \left( R(s,a) + \gamma \langle \underline{\Phi P(s,a)}, [V^*_M]_\phi \rangle \right)$$

$$= \sum_{s \in \phi^{-1}(x)} P_x(s) \left( \underbrace{R(s,a) + \gamma \langle P(s,a), V^*_M \rangle}_{} \right).$$

$$\underbrace{\phantom{\sum_{s \in \phi^{-1}(x)} P_x(s) \left( R(s,a) + \gamma \langle P(s,a), V^*_M \rangle \right)}}_{Q^*_M(s,a).}$$

# Loss of $\pi^{\star}_{M_\phi M}$ : approx. Q*-irrelevance

- Approximate case: proof breaks as $Q_M$* not piece-wise constant

- Workaround: define a new model $M_\phi'$ over $S$
$$R'_\phi(s,a) = \mathbb{E}_{\tilde{s} \sim p_{\phi(s)}}[R(\tilde{s},a)], \qquad P'_\phi(s'|s,a) = \mathbb{E}_{\tilde{s} \sim p_{\phi(s)}}[P(s'|\tilde{s},a)].$$

- Can show: $M_\phi$ and $M_\phi'$ share the same $Q$* (up to lifting)

- $\left\| [Q^{\star}_{M_\phi}]_M - Q^{\star}_M \right\|_\infty = \left\| Q^{\star}_{M_\phi'} - Q^{\star}_M \right\|_\infty \leq \dfrac{1}{1-\gamma} \left\| \mathcal{T}_{M_\phi'} Q^{\star}_M - Q^{\star}_M \right\|_\infty$

$$|(\mathcal{T}_{M_\phi'} Q^{\star}_M)(s,a) - Q^{\star}_M(s,a)|$$

$$= |R'_\phi(s,a) + \gamma \langle P'_\phi(s,a), V^{\star}_M \rangle - Q^{\star}_M(s,a)|$$

$$= \left| \left( \sum_{\tilde{s}:\phi(\tilde{s})=\phi(s)} p_x(\tilde{s}) \left( R(\tilde{s},a) + \gamma \langle P(\tilde{s},a), V^{\star}_M \rangle \right) \right) - Q^{\star}_M(s,a) \right|$$

$$= \left| \sum_{\tilde{s}:\phi(\tilde{s})=\phi(s)} p_x(\tilde{s}) \left( Q^{\star}_M(\tilde{s},a) - Q^{\star}_M(s,a) \right) \right| \leq \left| \sum_{\tilde{s}:\phi(\tilde{s})=\phi(s)} p_x(\tilde{s})(2\epsilon_{Q^\star}) \right| = 2\epsilon_{Q^\star}.$$

# Loss of $\pi^\star_{M_{\phi M}}$ :  approx. Q*-irrelevance

- Lesson: with Q*-irrelevance, the $\max_\pi \|V_M^\pi - V_{\widehat{M}}^\pi\|_\infty$ approach is not available; $\|Q_M^\star - Q_{\widehat{M}}^\star\|$ is the only choice

- If $\phi$ does not respect transition/reward, our analysis does not have to either!

# Recap

- **Theorem 2.** *(1) If $\phi$ is an $(\epsilon_R, \epsilon_P)$-approximate model-irrelevant abstraction, then $\phi$ is also an approximate $Q^\star$-irrelevant abstraction with approximation error $\epsilon_{Q^\star} = \frac{\epsilon_R}{1-\gamma} + \frac{\gamma \epsilon_P R_{\max}}{2(1-\gamma)^2}$.*
  *(2) If $\phi$ is an $\epsilon_{Q^\star}$-approximate $Q^\star$-irrelevant abstraction, then $\phi$ is also an approximate $\pi^\star$-irrelevant abstraction with approximation error $\epsilon_{\pi^\star} = 2\epsilon_{Q^\star}/(1-\gamma)$.*

- Given weighting distributions $\{p_x\}$, define $M_\phi = (S_\phi, A, P_\phi, R_\phi, \gamma)$

  $R_\phi(x, a) = \Sigma_{s \in \phi^{-1}(x)}\ p_x(s)\, R(s, a),\ \ P_\phi(x, a) = \Sigma_{s \in \phi^{-1}(x)}\ p_x(s)\, \Phi\, P(s, a).$

- How lossy is it to plan in $M_\phi$ and lift back to $M$?
  - If approx. bisimulation, use "$\max_\pi \|V_M^\pi - V_{\widehat{M}}^\pi\|_\infty$" type analysis

    $$\left\| V_M^\star - V_M^{[\pi^\star_{M_\phi}]_M} \right\|_\infty \leq \frac{2\epsilon_R}{1-\gamma} + \frac{\gamma \epsilon_P R_{\max}}{(1-\gamma)^2}$$

  - If approx. Q*-irrelevance, use "$\|Q_M^\star - Q_{\widehat{M}}^\star\|$" type analysis

    $$\left\| V_M^\star - V_M^{[\pi^\star_{M_\phi}]_M} \right\|_\infty \leq \frac{2\epsilon_{Q^\star}}{(1-\gamma)^2}$$

# Compare abstract model
# w/ bisimulation vs w/ Q*-irrelevance

Both guarantee optimality (exact case), but in different ways

- Consider value iteration (VI) in true model vs abstract model

- Bisimulation: every step of abstract VI resembles that step in true VI, throughout all iterations, b/c $\forall f : \phi(\mathcal{S}) \to \mathbb{R}, \ \mathcal{T}[f]_M = [\mathcal{T}_{M_\phi} f]_M$

- Q*-irrelevance: abstract VI initially behaves crazily. It only starts to resemble true VI when the function is close to $Q_M{}^*$

  - This is a circular argument

    $\mathcal{T}Q_M^\star = [\mathcal{T}_{M_\phi}[Q_M^\star]_\phi]_M$

  - Secret is stability—contraction of abstract Bellman update

  - Abstract Bellman update is a special case of projected Bellman update, and in general stability is not guaranteed. In that case, "Q*-irrelevance" alone is not enough to guarantee optimality

# The learning setting

- Given: $D = \{D_{s,a}\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ and $\phi$

- Algorithm: CE after processing data w/ $\phi$

- Shouldn't assume $|D_{s,a}|$ is the same for all $(s, a)$

  - … as we want to handle $|D| << |S|$

  - What should appear in the bound to describe sample size?

$$n_\phi(D) := \min_{x \in \mathcal{S}_\phi, a \in \mathcal{A}} |D_{x,a}|, \quad \text{where} \quad D_{x,a} := \bigcup_{s \in \phi^{-1}(x)} D_{s,a}.$$

  - At the mercy of data to be exploratory

# The learning setting

- Analysis varies according to whether $\phi$ is (approx.) bisimulation or Q*-irrelevant and the style ( $\max_\pi \|V_M^\pi - V_{\widehat{M}}^\pi\|_\infty$ vs $\|Q_M^\star - Q_{\widehat{M}}^\star\|$ )
- Will show analysis of Q*-irrelevance (can only use "$\|Q_M^\star - Q_{\widehat{M}}^\star\|$")
- Let $\widehat{M}_\phi$ be the estimated model
- Let $M_\phi$ be an abstract model w/ weighting distributions $p_x(s) \propto |D_{s,a}|$
- $M_\phi$ is the "expected model" of $\widehat{M}_\phi$
- 
$$\left\|Q_M^\star - [Q_{\widehat{M}_\phi}^\star]_M\right\|_\infty \leq \underbrace{\left\|Q_M^\star - [Q_{M_\phi}^\star]_M\right\|_\infty} + \underbrace{\left\|[Q_{M_\phi}^\star]_M - [Q_{\widehat{M}_\phi}^\star]_M\right\|_\infty}$$

*Approximation error*
- "Bias", informally
- Doesn't vanish with more data
- Smaller with a **finer** $\phi$
  (not w/ bisimulation; we will see why…)

*Estimation error*
- "Variance", informally
- Goes to 0 w/ infinite data
- Smaller with a **coarser** $\phi$

$$\left\| Q_M^\star - [Q_{\widehat{M}_\phi}^\star]_M \right\|_\infty \le \underbrace{\left\| Q_M^\star - [Q_{M_\phi}^\star]_M \right\|_\infty}_{\text{already handled}} + \underbrace{\left\| [Q_{M_\phi}^\star]_M - [Q_{\widehat{M}_\phi}^\star]_M \right\|_\infty}_{\text{to be analyzed}}$$

- Reusing the analysis for $\| Q_M^\star - Q_{\widehat{M}}^\star \|$
- Challenge: data is not generated from $M_\phi$
- Leverage the fact that Hoeffding can be applied to r.v.'s with non-identical distributions

$$\left\| [Q_{M_\phi}^\star]_M - [Q_{\widehat{M}_\phi}^\star]_M \right\|_\infty = \left\| Q_{M_\phi}^\star - Q_{\widehat{M}_\phi}^\star \right\|_\infty$$

$$\le \frac{1}{1-\gamma} \left\| Q_{M_\phi}^\star - \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^\star \right\|_\infty = \frac{1}{1-\gamma} \left\| \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^\star - \mathcal{T}_{M_\phi} Q_{M_\phi}^\star \right\|_\infty$$

$$|(\mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^\star)(x,a) - (\mathcal{T}_{M_\phi} Q_{M_\phi}^\star)(x,a)|$$

$$= |\widehat{R}_\phi(x,a) + \gamma \langle \widehat{P}_\phi(x,a), V_{M_\phi}^\star \rangle - R_\phi(x,a) - \gamma \langle P_\phi(x,a), V_{M_\phi}^\star \rangle|$$

$$= \left| \frac{1}{|D_{x,a}|} \sum_{s \in \phi^{-1}(x)} \sum_{(r,s') \in D_{s,a}} \left( r + \gamma V_{M_\phi}^\star(\phi(s')) - R(s,a) - \gamma \langle P(s,a), [V_{M_\phi}^\star]_M \rangle \right) \right|$$