

473 **A List of Notation**

Symbol	Type	Explanation
\perp		undefined
$\llbracket \text{bool} \rrbracket$	$\in \{0,1\}$	=1 if bool=True, =0 if bool=False
d	$\in \mathbb{N}$	number of states
k	$\in \mathbb{N}$	number of actions
i, j	$\in \mathbb{N}$	time index/step
$\{i:j\}$	$\subset \mathbb{Z}$	Set of integers from i to j (empty if $j < i$)
s, s', \dots, s^i	$\in \{1:d\}$	state at time step $1, 2, \dots, i$
a, a', \dots, a^i	$\in \{0:k-1\}$	action at time step $1, 2, \dots, i$
b, b', \dots, b^i	$\in \{0:k-1\}$	alternative action at time step $1, 2, \dots, i$
$a^{:i}$	$:= aa' \dots a^i$	sequence of i actions
$a^{<i}$	$:= aa' \dots a^{i-1}$	sequence of $i-1$ actions
\dot{s}, \ddot{s}	$\in \{1:d\}$	parts of state, usually $s = (\dot{s}, \ddot{s})$
ε	> 0	small number > 0
$p(\dots)$	$\in [0;1]$	(conditional) probability distribution over states and actions
$\pi(a s)$	$\in [0;1]$	policy. Probability of action a in state s
M^a, W^a	$\in [0;1]^{d \times d}$	transition-policy tensor $M_{ss'}^a = p(s' sa) \cdot \pi(a s)$ for each action a , $W = q$
B^a	$\in [0;1]^{d \times d}$	inverse 1-step model $B_{ss'}^a = p(a ss')$ for each action a
$B_{ss''}^{a++}$	$\in [0;1]$	3-step first-action inverse model $p(a ss''')$
J, K, Δ	$\in \mathbb{R}^{d \times d}$	action-independent $d \times d$ “transition” matrices
$\overset{+}{+}$	$\cdot^n \rightarrow \cdot$	index summation, e.g. $M_{s+}^+ = \sum_{a,s'} M_{ss'}^a$
\cdot	$(\cdot, \cdot) \rightarrow \cdot$	matrix multiplication: $[AB]_{ss''} = \sum_{s'} A_{ss'} B_{s's''}$
\odot	$(\cdot, \cdot) \rightarrow \cdot$	element-wise multiplication of matrix elements: $[A \odot B]_{ss'} = A_{ss'} B_{ss'}$
\oslash	$(\cdot, \cdot) \rightarrow \cdot$	element-wise division of matrix elements: $[A \oslash B]_{ss'} = A_{ss'} / B_{ss'}$
\otimes	$(\cdot, \cdot) \rightarrow \cdot$	tensor product: $[\dot{M} \otimes \ddot{M}]_{ss'} := \dot{M}_{\dot{s}\dot{s}'} \ddot{M}_{\ddot{s}\ddot{s}'}$ with $s = (\dot{s}, \ddot{s})$ and $s' = (\dot{s}', \ddot{s}')$

474 **B Characterizing M and W for which EqIM(1) holds**

$$M^a \oslash M^+ = W^a \oslash W^+ \iff W^a = M^a \odot J \quad \text{with} \quad J := W^+ \oslash M^+$$

475 That is, J is independent of a . Phrased differently

$$\text{For any } M \text{ and } W, \text{ EqIM(1) is satisfied iff } W^a \oslash M^a \text{ is independent } a. \quad (14)$$

476 For a given M , this allows to determine all W consistent with EqIM(1), by just multiplying with any
477 a -independent $J \geq 0$. Not all J though lead to W consistent with (7). In order to also satisfy (7), J
478 needs to be restricted as follows: With $\Delta_{ss'} := J_{ss'} - 1$, (7) becomes

$$0 \stackrel{!}{=} W_{s+}^a - M_{s+}^a = \sum_{s'} M_{ss'}^a (\Delta_{ss'} + 1) - M_{s+}^a = \sum_{s'} M_{ss'}^a \Delta_{ss'} \quad (15)$$

479 For each fixed s , these are k homogenous linear equations (one for each a) in d variables. Given M ,
480 all and only the W consistent with EqIM(1) and (7) can be obtained via $W^a = M^a \odot (1 + \Delta)$ with Δ
481 satisfying $M_s^+ \Delta_s = 0$.

482 As a special case, $\Delta = 0$ necessarily if and only if the rank of M_s^+ is $\geq d$ for every s . This gives the
483 precise conditions as stated in Proposition 1 under which (i) is true. We will next show that EqIM(2)
484 removes this limitation.

485 **C Characterizing M and W for which EqIM(1) and EqIM(2+) hold**

486 From Appendix B we know that the most general Ansatz for W^a satisfying EqIM(1) is $M^a \odot (1 + \Delta)$.
 487 Plugging this into (31) and expanding in Δ , we get

$$0 = M^a M^+ \odot (M^+)^2 - M^a M^+ \odot (M^+)^2 \\
+ M^a M^+ \odot [M^+ (M^+ \odot \Delta) + (M^+ \odot \Delta) \odot M^+] - [(M^a \odot \Delta) M^+ M^a (M^+ \odot \Delta)] \odot (M^+)^2 \\
+ M^a M^+ \odot (M^+ \odot \Delta)^2 - (M^a \odot \Delta) (M^+ \odot \Delta) \odot (M^+)^2$$

488 This is a collection of quadratic equations in Δ . The Δ -independent first line is 0. We can write this
 489 in canonical form:

$$\Sigma_{kl} A_{ss'',kl}^a \Delta_{kl} = R_{kl}^a(\Delta) \quad \text{with} \quad (16) \\
A_{ss'',kl}^a := (\Sigma_{s'} M_{ss'}^a M_{s's''}^+) (M_{sk}^+ M_{ks''}^+ \delta_{ls''} + M_{sl}^+ M_{ls''}^+ \delta_{sk} - M_{sk}^a M_{ks''}^+ \delta_{ls''} - M_{sl}^a M_{ls''}^+ \delta_{sk}) \\
R^a(\Delta) := (M^a \odot \Delta) (M^+ \odot \Delta) \odot (M^+)^2 - M^a M^+ \odot (M^+ \odot \Delta)^2$$

490 Let us consider A^a as a $d^2 \times d^2$ matrix for each a , Δ as a vector of length d^2 , and (wrongly) presume
 491 $R^a \equiv 0$ at first. A^a is a sum of 4 terms. The second and fourth terms are block-diagonal matrices
 492 (d blocks of size $d \times d$ in the diagonal) due to the δ_{sk} . The first and third terms are scrambled
 493 block-diagonal matrices due to the $\delta_{ls''}$, or more precisely, consist of $d \times d$ blocks, each block being
 494 a $d \times d$ diagonal matrix. If M^a has full rank, each of the four terms has full rank d^2 , but A^a itself
 495 can have lower rank, 0-eigenvalues due to some cancellations. Random M apparently achieves the
 496 highest rank, but even then, A^a itself has only rank $d(d-1)$.

497 Actually, $A^a \Delta = 0$ is required to hold for all a , so the rank of A as a $kd^2 \times d^2$ matrix may still be d^2 .
 498 But $A^+ \equiv 0$ for $k=2$ implies $A^0 = -A^1$, hence the rank is still at most $d(d-1)$. $k > 2$ may rectify
 499 this, but there is an alternative, which works for all a : Δ also needs to satisfy (15), which can be
 500 rewritten as

$$\sum_{kl} C_{s,kl}^a \Delta_{kl} = 0 \quad \text{with} \quad C_{s,kl}^a := M_{sl}^a \delta_{sk} \quad (17)$$

501 These give another kd constraints, and apparently often d new ones from random M . If we combine
 502 $A' := \begin{pmatrix} A \\ C \end{pmatrix}$, this implies that A' has often rank d^2 , so $A' \Delta = 0$ can only be satisfied for $\Delta = 0$. For
 503 $k=2$, $A^+ = 0$, so inclusion of either A^0 or A^1 in A' would suffice, but C^0 and C^1 are potentially
 504 independent, so both have to be included.

505 Let us now return to the real case of $R^a \neq 0$ for full random M , hence full-rank A' . With $R' := \begin{pmatrix} R \\ 0 \end{pmatrix}$,
 506 we need to solve $A' \Delta = R'$. Note that $R' = R'(\Delta)$ is not a constant, but a (homogenous) quadratic
 507 function of Δ itself. Consider any $\Delta = \Theta(\varepsilon)$, then $A' \Delta = \Theta(\varepsilon)$ while $R'(\Delta) = \Theta(\varepsilon^2)$, which is a
 508 contradiction for sufficiently small ε (this argument can be made rigorous). This implies that no Δ
 509 with $0 < \|\Delta\| < \varepsilon$ can satisfy $A' \Delta = R'(\Delta)$. In conclusion,

510 **Proposition 3 (Random M and full-rank A')**

511 *If A' has full rank and W is close to M , then EqIM(1) and EqIM(2) imply $W = M$*
 512 *Empirically A' has full rank for random M*

513 which of course implies EqIM(i) $\forall i$ and also (iv). Globally, i.e. if W is not close to M , these
 514 implications may not hold.

515 We have yet to establish sufficient conditions which M^a lead to full-rank A' . Empirically, this has
 516 been true for random M^a , so should hold almost surely if M are sampled uniformly. One might
 517 conjecture that full-rank M^a are sufficient, but this is not the case. For instance, if M^a is independent
 518 a , then $A' \equiv 0$.

519 **Zero A and R for full-rank \dot{M}^a .** We finally we note that A and R can have low rank, indeed $A \equiv$
 520 $0 \equiv R$ even for a -dependent full-rank M^a : Consider the example \dot{M}^a from (21) or its generalization
 521 (26): First, if for two matrices M^a and $M^{a'}$ only one s' (depending on s and s') contributes to
 522 the sum in $M^a M^{a'}$ then $(M^a \odot J)(M^{a'} \odot J) = M^a M^{a'} \odot K$ for some K . This makes (18) valid for
 523 $M^a := \dot{M}^a$ and $W^a := \dot{M}^a \odot J$ for any J , since for $aa' \neq bb'$ both sides are 0 by construction of \dot{M}^a
 524 (the $\odot K$ does nothing to it), and are trivially equal for $aa' = bb'$. By summing over $a'bb'$, also (31) is
 525 valid for any J , hence of course also for $J = 1 + \Delta$ for any Δ . Since (16) is equivalent to (31), (16)
 526 holds for any Δ . This can only be true for $A \equiv 0$ and $R \equiv 0$. This degeneracy in itself does not violate
 527 (ii), since the probability constraints require $W = M$, as established earlier.

528 **D EqIM(1) \wedge EqIM(2+) $\not\rightarrow$ EqIM(3) for full low rank M ?**

529 The following numerical approach may lead to counter-examples with full support to (v) without
 530 any divisions by 0 ($M_{ss'}^+ > 0$ and $W_{ss'}^+ > 0 \forall ss'$). We now consider full M^a but of rank $r < d$. The
 531 most interesting case is where all M^a span the same row-space, i.e. $M^a = L^a \cdot R$, where L^a are
 532 $d \times r$ matrices and R is a $r \times d$ matrix. Recall $A' := \begin{pmatrix} A^i \\ C^i \end{pmatrix}$ with A^a and C^a defined in (16) and (17).
 533 Empirically, for $k=2$, the rank of A' typically is $\min\{d^2, (3r-1)d - r(r-1)\}$, never more, and only
 534 in degenerate cases less. Hence for $r=2$, A' is singular for $d \geq 5$. Hence for $d \geq 5$, there exist $\Delta \neq 0$
 535 with $A' \Delta = 0$,

536 For $\Delta_0 := \Delta = \Theta(\varepsilon)$, this is an approximate $\Theta(\varepsilon^2)$ solution of $A' \Delta = R'(\Delta)$. By iterating $\Delta \leftarrow$
 537 $\Delta_0 + A'^+ R'(\Delta)$, where A'^+ is the pseudo-inverse of A' , we get an $\Theta(\varepsilon^i)$ -approximation after $i-2$
 538 iterations. This should rapidly converge to an “exact” non-zero(!) solution $A' \Delta = R'(\Delta)$. This would
 539 show that (ii) can fail for full M . Generically, this solution also violates EqIM(3), i.e. also (vi) can
 540 fail. By this we mean, for randomly sampled L^a and R (for $a=r=2$ and $d \geq 56$) and performing the
 541 procedure above, EqIM(3) does not hold. There is a caveat with this argument, namely if R' is not in
 542 the range of A' , then this construction fails.

543 **E EqIM(1) does not imply EqIM(2) (\odot -version)**

544 We have already given a simple example that violates (v) in Section 3, but the example and method-
 545 ology provided here generalizes to (vi) and even larger i . We consider deterministic reversible
 546 forward dynamics for any policy $\pi(a|s) > 0 \forall as$. For simplicity we assume $k=2$ and uniform policy
 547 $\pi(a|s) = \frac{1}{2}$. We defer a discussion of $0/0$ to the end of the next Appendix.

548 We consider M^a and W^a that permute states. That is, $M_{ss'}^a := \llbracket s' = \pi^a(s) \rrbracket$ and $W_{ss'}^a := \llbracket s' = \sigma^a(s) \rrbracket$
 549 for some permutations $\pi^a, \sigma^a : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$. Strictly speaking, we should multiply this
 550 by $\pi(a|s) = \frac{1}{k}$, but this global factor plays no role here, so will be dropped everywhere. Matrix
 551 multiplication corresponds to permutation composition: $[M^a W^a]_{ss''} = \llbracket s'' = \sigma^a(\pi^a(s)) \rrbracket$. We denote
 552 example permutation (matrices) by $[\pi] = [\pi(1) \dots \pi(d)]$.

553 We now construct a counter-example for (v): For $d=4$, let $M^0 = W^0 = \text{Id} = [1234]$ be the identity
 554 matrix/permutation. Let $W^1 = [2341]$ be the cyclic permutation $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 1$, and $M^1 = [2143]$
 555 the cycle pair $1 \leftrightarrow 2$ and $3 \leftrightarrow 4$. We know from (14) that EqIM(1) holds iff $W^a \odot M^a$ is independent
 556 $a (= J)$ iff $W^a \odot M^a = W^b \odot M^b \forall a, b \in \{0, 1\}$ iff $W^a \odot M^b = M^a \odot W^b$. Case $a = b$ is trivial,
 557 so only $W^0 \odot M^1 = M^0 \odot W^1$ needs to be verified. Now $M^a \odot W^a$ of two permutations matrices
 558 is *not* a permutation matrix (unless $M^a = W^a$). It still a 0-1 matrix with at most one non-zero
 559 entry in each row and column. We can generalize the permutation notation to “sub-permutations”
 560 by defining $\pi(s) = \emptyset$ if row s is empty. For instance $M^1 \odot W^1 = [2\emptyset 4\emptyset]$. EqIM(1) holds, since
 561 $W^0 \odot M^1 = [\emptyset\emptyset\emptyset\emptyset] = M^0 \odot W^1$.

562 Similarly EqIM(2a) holds iff $W^a W^{a'} \odot M^a M^{a'}$ is independent a, a' iff

$$W^a W^{a'} \odot M^b M^{b'} = M^a M^{a'} \odot W^b W^{b'} \quad \forall a, a', b, b'. \quad (18)$$

563 But for $a = a' = 0$ and $b = b' = 1$ we have

$$(W^0)^2 \odot (M^1)^2 = [1234] \odot [1234] = [1234] \neq [\emptyset\emptyset\emptyset\emptyset] = [1234] \odot [3412] = (M^0)^2 \odot (W^1)^2$$

564 hence EqIM(1) does not necessarily imply EqIM(2). The advantage of formulation (18) over (8) is
 565 that matrix sums M^+ and W^+ are more complicated objects than the sub-permutation matrices (18).
 566 Like random matrices, permutation matrices, have full rank, but unlike random matrices they can
 567 violate (ii), (iv), and (vi).

568 **F EqIM(1a) $\wedge \dots \wedge$ EqIM(ia) do not imply EqIM(i+1) (\odot -version)**

569 Counting variables and equations made the possibility of violating (v) for $k < d$ plausible (cf. positive
 570 result for $k \geq d$). A similar counting argument indicates that (vi) and higher i analogues might actually
 571 hold. Unfortunately this is not the case. I.e. even providing inverse models for all action sequences up
 572 to length i is not sufficient to always uniquely determine the probability of longer action sequences.
 573 This is true even for deterministic reversible forward dynamics for any policy $\pi(a|s) > 0 \forall as$. As for
 574 $i = 1$, we assume $k=2$, $\pi(a|s) = \frac{1}{2}$, gloss over $0/0$, and don't normalize M and W .

575 For $i=2$, $M^0 := W^0 := \text{Id} = [123456]$ and $W^1 := [234561] =: \sigma$ (σ for ‘cycle’) and $M^1 := [231564] =: \pi$
576 can be shown to satisfy EqIM(1) and EqIM(2a) but violate EqIM(3). The calculations are not to
577 onerous, but lets consider directly the general i case: Consider even $d =: 2d'$ and identity and cycle
578 (pair)

$$M^0 = W^0 = \text{Id} = [1, 2, \dots, d-1, d],$$

$$W^1 = [2, 3, \dots, d, 1], \quad M^1 = [2, 3, \dots, d', 1, d'+2, \dots, d-1, d, d'+1]$$

579 EqIM(ia) holds iff $W^a W^{a'} \dots \odot M^a M^{a'} \dots = W^+ W^+ \dots \odot M^+ M^+ \dots$ is independent $aa' \dots$ iff

$$W^a W^{a'} \dots W^{a^i} \odot M^b M^{b'} \dots M^{b^i} = M^a M^{a'} \dots M^{a^i} \odot W^b W^{b'} \dots W^{b^i} \quad \forall aa' \dots a^i, bb' \dots b^i \quad (19)$$

580 (While this looks like k^{2i} matrix equations, by chaining, checking k^i pairs suffices, which is the
581 same number as in EqIM(ia)). Now $W^a W^{a'} \dots W^{a^i}$ consists of only two types of matrices, a
582 cycle for $W^1 = \sigma$ and identity W^0 . The $W^0 = \text{Id}$ can be eliminated, leading to $(W^1)^{a^+}$, where
583 $a^+ := a + a' + \dots + a^i$. Similarly $M^b M^{b'} \dots M^{b^i} = (M^1)^{b^+}$, etc. Hence we only need to verify

$$(W^1)^{a^+} \odot (M^1)^{b^+} = (M^1)^{a^+} \odot (W^1)^{b^+} \quad \text{for } 0 \leq a^+, b^+ \leq i \quad (20)$$

584

$$(W^1)^{a^+} = [a^+ + 1, a^+ + 2, \dots, d, 1, 2, \dots, a^+], \quad \text{while}$$

$$(M^1)^{b^+} = [b^+ + 1, \dots, d', 1, \dots, b^+, d' + 1 + b^+, \dots, d, d' + 1, \dots, d' + b^+]$$

585 hence $(W^1)^{a^+} \odot (M^1)^{b^+} = [\emptyset \dots \emptyset] = 0$ for $0 \leq a^+ \neq b^+ < d'$. For $a^+ = b^+$ both sides of (20) are equal
586 too. Hence if we choose $d' = i + 1$, (20) and hence EqIM(1)...EqIM(ia) are all satisfied. If we choose
587 $d' = i$, $a^+ = d'$, $b^+ = 0$, (20) reduces to

$$(W^1)^{d'} \odot (M^1)^0 = [d' + 1, \dots, d, 1, \dots, d'] \odot \text{Id} = 0, \quad \text{and}$$

$$(M^1)^{d'} \odot (W^1)^0 = \text{Id} \odot \text{Id} = \text{Id}$$

588 which are of course not equal. Hence EqIM(i) fails for $d' = i$. Summing over all $a' \dots a^{d'}$ and $b' \dots b^{d'}$,
589 and noting that all other terms are 0 or cancel, shows that EqIM($i+$) fails too. Together this shows
590 for $d' = i + 1$ that EqIM(1)...EqIM(ia) do not imply any version of EqIM($i + 1$).

591 Despite M^a having full rank, A and A' defined in Appendix C have very low rank, indicating
592 potentially many more consistent W .

593 A downside of this example is that it strictly only applies to the \odot -version (19). Many entries of
594 M^+ and W^+ and powers thereof are 0, so (8) contains many divisions by zero. We were not able to
595 extend this example by mixing in e.g. a uniform matrix as done in the first counter-example to (v).

596 Many real-world MDPs are sparse. Only a subset $G \subseteq S \times S$ of transitions $s \rightarrow s'$ is possible. For
597 $(s, s') \notin G$, $p(s' | sa) = 0 \forall a$, or formally $M_{ss'}^a = M_{ss'}^+ = 0$. In this case, no action causes $s \rightarrow s'$ and
598 $p(a | ss') = M_{ss'}^a / M_{ss'}^+$ being undefined is actually appropriate. So we could restrict (s, s') to G (and
599 analogously (s, \dots, s^i) and (ss^i) by chaining G) in the conditions and conclusions of the various
600 conjectures. It is then also natural to restrict the model class to $\mathcal{M} := \{M : M_{ss'}^+ > 0 \Leftrightarrow (s, s') \in G\}$.
601 For unknown G , the condition $M, W \in \mathcal{M}$ then becomes $M_{ss'}^+ > 0 \Leftrightarrow W_{ss'}^+ > 0$. Unfortunately
602 the above counter-example does not even satisfy this weaker condition, but the more complicated
603 example of Appendix G does.

604 G Non-Uniqueness of Inverse MDP Models for $i \geq 2$

605 In Appendices E/F we provided conjectured/unsatisfactory counter-examples to EqIM(1 : i) \Rightarrow
606 EqIM($i + 1$). Here we provide a fully satisfactory counter-example that avoids the ‘‘bad’’ 0/0.

607 **EqIM(1) and EqIM(2a) do not imply EqIM(3).** Consider two matrices \dot{M}^0 and \dot{M}^1 with
608 disjoint support, i.e. $\dot{M}^0 \odot \dot{M}^1 = 0$. In this case $\dot{M}^a \odot \dot{M}^+ \in \{0, 1, \perp\}^{d \times d}$ is a partial binary matrix
609 with entry undefined (\perp) wherever $\dot{M}^+ = 0$ but otherwise 0 wherever $\dot{M}^a = 0$ and 1 wherever $\dot{M}^a > 0$.
610 That is, it is insensitive to the actual (non-zero) values of \dot{M}^a . A simple such \dot{M} is $\dot{M}^0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ and
611 $\dot{M}^1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, ignoring normalization. For now we ignore ss' for which $\dot{M}_{ss'}^+ = 0$ and return to this
612 issue later.

613 We consider M^a and W^a that permute states. That is, $M_{ss'}^a := \llbracket s' = \pi^a(s) \rrbracket$ and $W_{ss'}^a := \llbracket s' = \sigma^a(s) \rrbracket$
614 for some permutations $\pi^a, \sigma^a : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$. Strictly speaking, we should multiply this by
615 e.g. $\pi^a(a|s) = \frac{1}{k}$, but this global factor plays no role here, so will be dropped everywhere. Matrix
616 multiplication corresponds to permutation composition: $[M^a W^a]_{ss''} = \llbracket s'' = \sigma^a(\pi^a(s)) \rrbracket$. We denote
617 example permutation (matrices) by $[\pi] = [\pi(1) \dots \pi(d)]$. Consider now

$$\begin{aligned} \dot{M}^0 &:= [456123] =: [\pi_0] & \implies & \dot{M}^0 \dot{M}^1 = [564312] & (21) \\ \dot{M}^1 &:= [231645] =: [\pi_1] & & \dot{M}^1 \dot{M}^0 = [645231] \\ & & & \dot{M}^1 \dot{M}^1 = [312564] \end{aligned}$$

618 No column contains the same number twice, hence this not only satisfies $\dot{M}^0 \circ \dot{M}^1 = 0$ but also

$$\dot{M}^a \dot{M}^{a'} \circ \dot{M}^b \dot{M}^{b'} = 0 \quad \text{unless } a=b \text{ and } a'=b' \quad (22)$$

619 That $6 \rightarrow 5 \rightarrow 4 \rightarrow 6$ is in reverse order to $1 \rightarrow 2 \rightarrow 3 \rightarrow 1$ is crucial for making \dot{M}^0 and \dot{M}^1 not commute.
620 Note that (22) remains valid if each 1-entry of \dot{M}^a is replaced by a different non-zero scalar, since
621 (22) is purely multiplicative. So if $\dot{W}^a = \dot{M}^a \circ J$ for some $J > 0$, then $\dot{W}^a \dot{W}^{a'} = \dot{M}^a \dot{M}^{a'} \circ K$ for
622 some $K > 0$. Let \dot{W}^a be such a matrix. Then $[\dot{W}^a \dot{W}^{a'} \circ \dot{W}^+ \dot{W}^+]_{\dot{s}\dot{s}''} = 1$ if $[\dot{M}^a \dot{M}^{a'}]_{\dot{s}\dot{s}''} > 0$ and 0
623 (or undefined) otherwise, i.e. is independent of the choice of J . So such $\dot{W} \neq \dot{M}$ satisfies EqIM(2a).
624 Unfortunately the probability constraints $W_{s+}^a = 1$ require $J_{s+}^a = 1$ when $M_{ss'}^+ > 0$, and hence $W = M$.
625 But the general idea is sound and can be made work as follows:

626 We split one state, e.g. $s = 6$ into two states $s = 6a$ and $s = 6b$. We leave the permutation structure
627 intact, except that all deterministic transitions into $s = 6$ are split into stochastic transitions to $s = 6a$
628 and $s = 6b$, and transitions from $6a$ and $6b$ will be to the same state as from original 6. Condition (22)
629 is still satisfied, so the above argument still goes through, but now we can choose different stochastic
630 transitions to $s = 6a$ and $s = 6b$ in W and M .

631 Finally, we have to show violation of EqIM(3). EqIM(ia) holds iff $W^a W^{a'} \dots \circ M^a M^{a'} \dots =$
632 $W^+ W^+ \dots \circ M^+ M^+ \dots$ is independent $aa' \dots$ iff

$$W^a W^{a'} \dots W^{a^i} \circ M^b M^{b'} \dots M^{b^i} = M^a M^{a'} \dots M^{a^i} \circ W^b W^{b'} \dots W^{b^i} \quad \forall aa' \dots a^i, bb' \dots b^i \quad (23)$$

633 (While this looks like k^{2i} matrix equations, by chaining, checking k^i pairs suffices, which is the same
634 number of equations as in EqIM(ia)).

635 It is easier to split every state into two states: $s := (\dot{s}, \ddot{s})$ with $\dot{s} \in \{1, \dots, 6\}$ as before and splitter
636 $\ddot{s} \in \{0, 1\}$. $M_{ss'}^a := \dot{M}_{\dot{s}\dot{s}'}^a \ddot{M}_{\ddot{s}\ddot{s}'}^a$. Note that \ddot{M} is flexible enough to expand each 1-entry in \dot{M}^a to a
637 different 2×2 (stochastic) matrix, while the 0-entries become $\begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$. This flexibility is important: \ddot{M}
638 independent a or independent \dot{s} would not work. Now let us write out

$$[M^a M^{a'} M^{a''}]_{ss'''} = \sum_{\dot{s}' \dot{s}''} \dot{M}_{\dot{s}\dot{s}'}^a \ddot{M}_{\ddot{s}' \dot{s}''}^{a'} \dot{M}_{\dot{s}'' \dot{s}'''}^{a''} \sum_{\ddot{s}' \ddot{s}''} \dot{M}_{\dot{s}\dot{s}'}^{a\dot{s}} \ddot{M}_{\ddot{s}' \dot{s}''}^{a' \dot{s}'} \ddot{M}_{\dot{s}'' \dot{s}'''}^{a'' \dot{s}''} \quad (24)$$

639 The crucial difference to the $i=2$ case (22) is that now there are difference permutation sequences
640 leading to the same permutation, for instance $\dot{M}^0 \dot{M}^0 \dot{M}^1 = \dot{M}^1 = \dot{M}^1 \dot{M}^0 \dot{M}^0$. Let us choose
641 $aa'a'' = 001$ and $\dot{s} = 1$, then only $\dot{s}' = \pi_0(\dot{s}) = 4$ and $\dot{s}'' = \pi_0(\dot{s}') = 1$ contribute to the sum and
642 $\dot{s}''' = \pi_1(\dot{s}'') = 2$. For this choice, (24) becomes $1 \cdot 1 \cdot 1 \cdot [\dot{M}^{01} \dot{M}^{04} \dot{M}^{11}]_{\dot{s}\dot{s}''}$. If we replace $aa'a''$ in
643 (24) by $bb'b''$ and then choose $bb'b'' = 100$ and again $\dot{s} = 1$, then only $\dot{s}' = \pi_1(\dot{s}) = 2$ and $\dot{s}'' = \pi_0(\dot{s}') = 5$
644 contribute and $\dot{s}''' = \pi_0(\dot{s}'') = 2$. For this choice, (24) becomes $1 \cdot 1 \cdot 1 \cdot [\dot{M}^{11} \dot{M}^{02} \dot{M}^{05}]_{\dot{s}\dot{s}''}$. We now
645 define $W_{ss'}^a := \dot{M}_{\dot{s}\dot{s}'}^a \ddot{W}_{\ddot{s}\ddot{s}'}^{a\dot{s}}$. Since \dot{M} remains the same, the same action and state sequences above
646 lead to the same result for W , just with \ddot{M} replaced by \ddot{W} . If we plug the four expressions into (23)
647 (for $i=3$) we get

$$\ddot{W}^{01} \ddot{W}^{04} \ddot{W}^{11} \circ \ddot{M}^{11} \ddot{M}^{02} \ddot{M}^{05} = \ddot{M}^{01} \ddot{M}^{04} \ddot{M}^{11} \circ \ddot{W}^{11} \ddot{W}^{02} \ddot{W}^{05}$$

648 Since this expressions involves 10 different 2×2 stochastic matrices, there are plenty of choices to
649 make both sides different. If we choose all 2×2 matrices to have full support, then by construction,
650 W and M have the same support, hence constitute a proper counter-example to EqIM(3). We now
651 extend this construction to $i > 2$.

652 **EqIM(1a) \wedge ... \wedge EqIM(ia) do not imply EqIM(i+1).** The construction in the previous para-
653 graph generalizes to $i > 2$: We need to find two permutations $\dot{M}^0 = \pi_0$ and $\dot{M}^1 = \pi_1$ such that for each
654 fixed $j \leq i$ all possible 2^j concatenations (products) of these permutation (matrices) differ in the sense
655 that no s is mapped to the same s^j (they have disjoint support). Since all $\dot{M}^a \dot{M}^{a'} \dots \dot{M}^{a^j} \in \{0,1\}$, we
656 can write this condition compactly as

$$\sum_{aa' \dots a^j} \dot{M}^a \dot{M}^{a'} \dots \dot{M}^{a^j} \in \{0,1\}^{d \times d}$$

657 By factoring the sum, this is equivalent to $(\dot{M}^+)^j \in \{0,1\}^{d \times d}$. Note that $[(\dot{M}^+)^j]_{ss^i}$ counts the
658 number of action sequences $aa' \dots a^j$ of length j that lead from s to s^i . For $j = i+1$, we want this
659 condition to be violated. So in order to disprove the implication we need to find two permutations
660 M^0 and M^1 such that

$$(\dot{M}^+)^j \in \{0,1\}^{d \times d} \quad \forall j \leq i \quad \text{but} \quad (\dot{M}^+)^{i+1} \notin \{0,1\}^{d \times d} \quad (25)$$

661 The rest of the argument is the same as for the $i=2$ case above: creating two versions M^a and W^a of
662 \dot{M}^a by spitting one or all states into two, and replacing the 1s by 2×2 different stochastic matrices.
663 As for the choice of \dot{M}^a , for $i=3$ we can choose 3-cycle and 5-cycle

$$\begin{aligned} \dot{M}^0 &= [6,7,8,9,10,11,12,13,14,15,1,2,3,4,5] \\ &= (1,6,11)(2,7,12)(3,8,13)(4,9,14)(5,10,15) \\ \dot{M}^1 &= [2,3,4,5,1,8,9,10,6,7,14,15,11,12,13] \\ &= (1,2,3,4,5)(6,8,10,7,9)(11,14,12,15,13) \end{aligned} \quad (26)$$

664 where we also provide the more conventional cycle notation in round brackets. Crucially the 5-cycles
665 have been chosen to not commute with the 3-cycles ($M^0 M^1 \neq M^1 M^0$). Conditions (25) can easily
666 be verified numerically. For higher i we need p cycles and q cycles, where p and q are relative
667 prime and sufficiently large. We need at least $d = p \cdot q \geq 2^i$, otherwise $\dot{M}^+ \notin \{0,1\}^{d \times d}$ by a simple
668 pigeon-hole argument. To prove $\text{EqIM}(1a) \wedge \dots \wedge \text{EqIM}(ia) \not\Rightarrow \text{EqIM}(i+1)$ in general for arbitrarily
669 large i , we need to invoke some group theory. All-together we have shown that

670 **Proposition 4 ((i)-(vi) can fail)** *EqIM(1a) \wedge ... \wedge EqIM(ia) do not necessarily imply EqIM(i+1) for*
671 *any i . This in turn implies that (i)-(vi) each can fail for some M .*

672 H Deterministic Cases

673 **Deterministic planning / reachability problem.** If we are only interested in finding *some* action
674 sequence $aa' \dots a^i$ that leads to s^i , the problem becomes easy: The only thing that matters is the
675 support of the various matrices, not the numerical values themselves. Since $B_{ss'}^a > 0$ iff $M_{ss'}^a > 0$
676 (either assuming $M_{ss'}^+ > 0$ or regarding $\perp > 0$ as False), and similarly for higher orders, we can replace
677 M^a by B^a in (iii), and get $B_{ss^{i+1}}^{aa' \dots a^i} > 0$ iff $[B^a B^{a'} \dots B^{a^i}]_{ss^{i+1}} > 0$. We could also replace M^a by
678 $G_{ss'}^a := \llbracket B_{ss'}^a > 0 \rrbracket$, then $[G^a G^{a'} \dots G^{a^i}]_{ss^{i+1}} > 0$ counts the number of paths of length i from s to s^{i+1}
679 via action sequence $aa' \dots a^i$, and hence determines whether s^{i+1} can be reached. Similarly $(G^+)^i > 0$
680 iff there is *some* action sequence that can reach s^{i+1} from s . An action a such that $G^a (G^+)^i > 0$ can
681 be chosen as the first action of such a sequence if it exists, and $a', a'' \dots$ can be found the same way by
682 recursion. So this deterministic planning/reachability problem has a “unique” solution, which can be
683 found in time $O(i \cdot d \cdot (d+k))$ (for fixed s and s^{i+1}).

684 **B is deterministic.** Assume $M_{ss'}^a / M_{ss'}^+ =: B_{ss'}^a \in \{0,1,\perp\}$. This is true if and only if M^a has
685 disjoint support for different a , i.e. iff $M^a \odot M^b = 0 \quad \forall a \neq b$. This in turn means that $B_{ss'}^a = \llbracket W_{ss'}^a > 0 \rrbracket$
686 for any and only those W with same support as M , and hence also $W^a \odot W^b = 0 \quad \forall a \neq b$, which is
687 another failure case of (i). Here we have included the case where *no* action leads from s to s' , in which
688 case $W_{ss'}^+ = 0$ and B^a is undefined (\perp). This readily extends to higher orders: If $B^{aa' \dots} \in \{0,1,\perp\}$,
689 then $B^{aa' \dots} = \llbracket W^a W^{a'} \dots \odot (W^+)^i > 0 \rrbracket$ iff $W^a W^{a'} \dots$ has the same support as $M^a M^{a'} \dots$ and

$$W^a W^{a'} \dots W^{a^i} \odot W^b W^{b'} \dots W^{b^i} = 0 \quad \forall aa' \dots a^i \neq bb' \dots b^i \quad (27)$$

690 Note that $W^a \odot W^b = 0$ does not necessarily imply (27), e.g. for $W^0 = \frac{1}{2} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ and $W^1 = \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$,
691 $(W^0)^2 = (W^1)^2$. In Appendices G&E&F we construct W such that (27) holds for larger i .

692 I Applications

693 Consider an agent who has control over \dot{s} but not over \ddot{s} . For instance a robot equipped with a
 694 camera can control its position and orientation, but not the shape and color of objects in its path.
 695 The forward model $p(s'|as)$ essentially involves modelling the whole observable world. The inverse
 696 model $p(a|ss')$ on the other hand can ignore inputs that the agent has no control over. Of course
 697 in practice, s does not come neatly separated into \dot{s} and \ddot{s} , so a (say) deep neural network still has
 698 to learn the controllable features, but neither needs to learn nor predict the uncontrollable features
 699 (under the factorization assumptions described in Section 3, now in feature space).

700 If the goal is to navigate from s to s^i in i time steps, and open-loop control suffices as e.g. in
 701 (near)-deterministic problems [EMK⁺22], then action sequences for which $p(aa' \dots a^{i-1} | ss^i)$ is large
 702 are the most likely that caused the transition to s^i , hence these sequences are promising candidates
 703 for macro actions (temporally extended actions, options) in Reinforcement Learning [SP02, Pre00].

704 Since the action space is typically much smaller than the state space (the former often finite, the
 705 latter often even infinite-dimensional), even learning $p(aa' \dots a^{i-1} | s \dots s^i)$ directly for all small i can
 706 be feasible and may be more efficient than learning the one-step forward model. A closed-loop
 707 alternative would be to learn only $p(a | s \dots s^i)$, find the likely first action a that caused the ultimate
 708 transition to s^i , then take action a , iterate, and store the resulting sequence as an option.

709 The required sample complexity to learn inverse MDP models for larger i directly from data may
 710 grow exponentially in i , which is why inferring i -step inverse models from 1-step and 2-step inverse
 711 models would be useful. The fact that this problem borders NP-hardness probably prevents even
 712 powerful transformer models to finding the structure in $p(aa' \dots a^{i-1} | s \dots s^i)$ by themselves.

713 J Systems of Quadratic Matrix Equations

714 A System of Polynomial Equations (SPE) is a set of multivariate polynomial equations
 715 $\text{Poly}_j(x, y, z, \dots) = 0$ over \mathbb{R} in n variables $x, y, z, u, v, w, \dots \in \mathbb{R}$ for $j \in \{1 : m\}$. This class is NP-
 716 hard (via a simple reduction from 1in3SAT, see Section K). We can recursively replace each product
 717 xy (sum $bu + cv$) in the polynomials by a new variable z (w) and add “polynomial” equation $z = xy$
 718 ($w = bu + cv$). This results in SPEs consisting of only linear equations with a single $+$ ($bu + cv = w$)
 719 and quadratic equations without any $+$ ($xy = z$), which are still (even with all $a = b = 1$ and $x = y = z$)
 720 NP-hard. We call them Simple Systems of Quadratic Equations (Simple SQE). For the reduction pro-
 721 cess to actually work we need one further dummy variable and equation $q = 1$ (to reduce $bu + cv = w$).
 722 Alternatively, with some extra work, we can reduce any SPE into a Simple SQE asking for a *non-zero*
 723 solution. We will pursue the latter, since this is closer to our interest (SQE (16) with solution $\Delta \neq 0$).
 724 We can even merge the linear and quadratic equations into a single form $xy = bu + cv$ by choosing
 725 $b = 1$ and $c = 0$ (replacing xy by w and adding $xy = 0 \cdot u + 1 \cdot w$).

726 We define a System of Polynomial/Quadratic Matrix Equations (SPME/SQME) as a set of m
 727 multivariate (quadratic) polynomials $\text{Poly}_j(\Delta, \Gamma, \dots | A, B, C, \dots) = 0$ in the (unknown) matrix variables
 728 Δ, Γ, \dots and the (given) matrix constants (“coefficients”) A, B, C, \dots . Alternatively, Poly_j might be
 729 viewed as generalized polynomials over a *non-commutative* matrix ring in the unknowns only. In any
 730 case, note that

$$A \cdot \Delta \cdot A' \cdot \Delta \cdot A'' + B \cdot \Delta \cdot B' + C \neq (A \cdot A' \cdot A'') \cdot \Delta^2 + (B \cdot B') \cdot \Delta + C$$

731 By writing out all matrix operations in terms of their scalar operations, SPME is of course a sub-class
 732 of SPE. SPE is also a sub-class of SPME (choose all matrices to be 1×1 matrices), which implies
 733 SPME is NP-hard. But we are interested in NP-hard *small* subclasses of SPME, so will construct
 734 a more economical embedding: Assume we have a Simple SQE with n variables x, y, z, u, v, \dots . We
 735 place them into $d \times d$ matrix Δ ($d \geq \sqrt{n}$) introducing dummy variables for the remaining entries. We
 736 can extract variable $w = \Delta_{ss'}$ via $w = e^{s\top} \cdot \Delta \cdot e^{s'}$, where e^s is basis vector ($d \times 1$ -matrix) (e^s) _{$s'1$} = $\delta_{ss'}$.
 737 If we replace all variables in the Simple SQE expressions $xy = au + bv$ by such expressions, we get a
 738 Simple SQME with Poly_j equations of the form (dropping \cdot as usual)

$$a^j \Delta A'^j \Delta a''^j = b^j \Delta b'^j + c^j \Delta c'^j \quad \forall j \quad (28)$$

739 While these are scalar equations, since the outer matrices are $1 \times d$ on the left and $d \times 1$ on the right,
 740 technically they are matrix equations. We could pad all involved matrices, including the outer ones,

741 with zeros to square $\mathbb{R}^{d \times d}$ matrices of the same size (for sufficiently large d , and only polynomial
742 overhead).

743 We can reduce (28) to just one equation at the cost of making the equations more complicated as
744 follows: Write each equation $\text{Poly}_j = 0$ in the form $e^s \cdot \text{Poly}_j \cdot e^{s' \top} = 0$, with a different (s, s') -pair for
745 each j . These are now “proper” matrix equations, but with all entries identically 0 except entry (s, s')
746 being Poly_j . This allows us to sum all equations without conflating them into one (complex) matrix
747 equations

$$\sum_j A^j \Delta A'^j \Delta A''^j = \sum_j B^j \Delta B'^j + C^j \Delta C'^j \quad (29)$$

748 Another way to combine (28) into one equation is by putting all M^j for all j into one block-diagonal
749 matrix $\tilde{M} := \text{Diag}(M^1, \dots, M^m)$ for $M \in \{a, A', a'', b, b', c, c', \Delta\}$. For $\tilde{\Delta}$ we need to ensure that indeed
750 all blocks $\Delta^j = \Delta$ are equal. This can be done via $\tilde{\Pi}^\top \tilde{\Delta} \tilde{\Pi} = \Delta$ for some cyclic block permutation $\tilde{\Pi}$.
751 We further need to ensure that the off-diagonal blocks of $\tilde{\Delta}$ are zero. We can zero each block with
752 one equation, but it seems impossible to zero all with a bounded number of Simple QMEs. We can
753 modify the decision problem to decide whether specific sparse solutions $\tilde{\Delta}$ exist. Formally, we can
754 introduce element-wise multiplication \odot and allow one equation of the form $\tilde{B} \odot \tilde{\Delta} = 0$ with \tilde{B} being
755 0/1 on the on/off-diagonal blocks. This leads to a Simple SQME with \odot in 3 equations (dropping the
756 \sim)

$$A \Delta A' \Delta A'' = B \Delta B' + C \Delta C', \quad \Pi^\top \Delta \Pi = \Delta, \quad B \odot \Delta = 0 \quad (30)$$

757 **Proposition 5 (NP-hardness of Simple SQME)** *Systems of Polynomial Equations (SPE) can be*
758 *polynomially reduced to Simple Systems of Quadratic Matrix Equations (Simple SQME) (28). The*
759 *number of equations can be reduced to 1 at the expense of making the equations complex (29), or to*
760 *2 by asking for sparse solutions or by enforcing sparsity via $B \odot \Delta = 0$ (30). Since SPE are NP-hard,*
761 *deciding the existence of non-zero solutions for all three SQME versions is also NP-hard.*

762 An NP-hardness proof for a Simple SQME with \odot with 3 equations via reduction from lin3SAT
763 that looks much closer to the desired form (32) or (34) is given in Section K. By a similar reduction,
764 encoding all n variables and their complement in the diagonal of $\Delta = \text{Diag}(x, \bar{x}, y, \bar{y}, \dots)$, one can also
765 show that solvability of

$$\Delta^2 = \Delta, \quad A \Delta 1 = 1, \quad \text{Id} \odot \Delta = \Delta, \quad \text{with } A \in \{0, 1\}^{m \times 2n}$$

766 is NP-complete (1 is the all-1 vector, sparse A with 2 or 3 ones in each row suffice), but not all SPE
767 can be reduced to this form.

768 **Open Problem 6 (Are Bounded SPME NP-hard?)** *Are Systems of Polynomial Matrix Equations*
769 *(without \odot) of bounded structural complexity NP-hard? Bounded means, only the definitions of the*
770 *constant matrices scale with $d \times d$, but the polynomial degrees, number of equations, and number of*
771 *matrix operations are bounded.*

772 K Computational Complexity

773 Maybe even just characterizing all M for which EqIM(1) and EqIM(2) uniquely determine W is
774 hopeless, not to speak of finding some or all W in case not. More formally, we can ask the question
775 of whether there exists an efficient algorithm that can decide whether EqIM(i) has a unique solution.
776 We provide some weak preliminary evidence, why this problem may be NP-hard. Appendix M
777 contains fully self-contained a few versions of this open problem in their simplest instantiation and
778 most elegant form.

779 **Decidability and computability.** EqIM(2) converted to (23) and (7), or (31) or (32) below form a
780 System of Quadratic Equations (SQE). The constraint $W \neq M$ can also be expressed as a quadratic
781 equation (see below). As such, the existence and uniqueness of solutions is formally decidable
782 by computing a Gröbner basis [Stu02], and (some) solutions can be found by cylindrical algebraic
783 decomposition in (double) exponential time. ε -approximate solutions can of course be found by
784 exponential brute-force search through all W on a finite ε' -grid, and verified in polynomial time.

785 **Complexity considerations.** 3SAT is NP complete. A CNF formula in n boolean variables can
786 easily be converted to a System of Quadratic Equations (SQE). Therefore SQE is also NP hard.
787 EqIM(2+) explicitly written in quadratic form

$$M^a M^+ \odot (W^+)^2 - W^a W^+ \odot (M^+)^2 = 0 \quad (31)$$

788 constitutes an SQE in W given M , also if we include linear EqIM(1) and probability constraints
789 (7). Non-negativity of W can be enforced with (slack) variables $(Y_{ss'}^a)^2 = W_{ss'}^a$. (Similarly (16)
790 plus constraints (15) constitute an SQE in Δ .) To reduce the uniqueness question to a solvability
791 problem we need to avoid the trivial solution $W \equiv M$, e.g. by introducing further (slack) variables
792 $t \in \mathbb{R}$ and $\Gamma_{ss'}^a := (W_{ss'}^a - M_{ss'}^a)^2$ and constraint $t \cdot \Gamma_{++}^+ = 1$. Due to the minus sign in (31), this cannot
793 be converted to a convex (optimization) problem. The choice of M gives significant freedom in
794 creating SQE problems, even if only considering permutation matrices $M^a \in \{0,1\}^{d \times d}$. If one could
795 show that every SQE can be represented as (31) [plus $W \neq M$ constraint] for a suitable choice of M ,
796 this would imply that proving the existence of $W \neq M$ satisfying (31) is NP hard. This in turn would
797 imply that computing (any) $p(a|ss''')$ from $p(a|ss')$ and $p(a|ss'')$ is NP hard. On the other hand,
798 matrix multiplication $W^a W^b$ is a very specific quadratic form, which may not be flexible enough to
799 incorporate every SQE within (31).

800 We could not find any work on NP-hardness of Systems of Polynomial Matrix Equations (SPME).
801 There is work on the NP-hardness of tensor problems [HL13], but this refers to the design tensors, e.g.
802 $\sum_{jk} A_i^{jk} x_j x_k + \sum_j B_i^j x_j + C_i = 0 \forall i$, but the unknowns are always treated as scalars or vectors. Of
803 course $[X \cdot Y]_{ik} = \sum_{abcd} A_{ik}^{abcd} X_{ab} Y_{cd}$, but $A_{ik}^{abcd} = \delta_{ai} \delta_{dk} \delta_{bc}$ is a very special fixed tensor (actually
804 of low tensor rank d) with no flexibility of encoding NP-hard problems therein.

805 That inference in Bayesian networks is NP-complete [KF09] does not help us either for two reasons:
806 First, in our problem the probability distribution over states and actions is only partially given. More
807 importantly, our network for $i=2$ has only 5 nodes (s, a, s', a', s'') , while the NP-hardness proofs we
808 are aware of require large networks. Even for fixed $i > 2$, it is not obvious how to encode NP-hard
809 problems into EqIM(i), due to the severe structural constraints in EqIM(i) compared to a general
810 network with $2i+3$ nodes. It is not clear how to exploit the fact that our (few) state nodes are large.

811 SQE are polynomially equivalent to Systems of Quadratic Matrix Equations (SQME), which may be
812 the reason complexity theorists have ignored the latter. We suspect but do not know whether SQME
813 of *bounded* structural complexity (only the definitions of the constant matrices scale with $d \times d$) is
814 NP-hard (Open Problem 6). If we allow sparse encoding of SQE variables in W , i.e. we allow one
815 equation involving \odot of the form $B \odot W = 0$ with boolean matrix B , then bounded SQME becomes
816 NP-hard. See Appendix J for details.

817 Below we directly reduce 1in3SAT to a Bounded-SQME with \odot that resembles our problem as close
818 as we were able to make it.

819 **An NP-hard matrix problem.** From EqIM(1) we know that $W^a = B^a \odot W^+$. Plugging this into
820 EqIM(2a) gives

$$B^{aa'} \odot (W^+ \cdot W^+) = (B^a \odot W^+)(B^{a'} \odot W^+) \quad \text{with constraints} \quad [B^a \odot W^+]_{s+} = \pi(a|s) \quad (32)$$

821 This set of equations is purely in terms of what is given (B^a and $B^{aa'}$) and only involves unknowns
822 W^+ without reference to W^a . See Appendix L for some further simplification and discussion. We
823 will show:

824 **Proposition 7 (An NP-hard matrix problem)** *Given A, B, C, Π , deciding whether the following*
825 *quadratic matrix problem has a solution in W is NP-hard:*

$$A \odot (W \cdot W) = (C \odot W)(C \odot W), \quad [B \odot W]_{s+} = 1, \quad \Pi \cdot W = W \quad (33)$$

826 This has some resemblance to (32). Since the boundary between P and NP is very fractal/subtle,
827 this in-itself may not imply much, but is more meant as a demonstration of how one may approach
828 proving NP-hardness of (32).

829 **Proof.** We reduce 1in3SAT, which is an NP-complete variant of 3SAT, where each clause must have
830 exactly one satisfying assignment, to (33). A 3CNF(n, m, g) formula is a boolean conjunction of m
831 clauses in n variables, where each clause $c_i = \ell_{i1} \vee \ell_{i2} \vee \ell_{i3}$ for $i \in \{1 : m\}$ is a 1-in-3 disjunction of 3

832 literals, and each literal is $\ell_{ia} = x_j$ or it's complement $\ell_{ia} = \neg x_j \equiv \bar{x}_j$, where $j = g(i, a)$ is the variable
833 index of clause i in position a .

834 We arithmetize the 3CNF expression in the standard way by replacing True $\rightsquigarrow 1$, False $\rightsquigarrow 0$, and $\vee \rightsquigarrow +$,
835 i.e. we ask whether the system of linear equations $\ell_{i1} + \ell_{i2} + \ell_{i3} = 1 \forall i$ has a solution in $x_j \in \{0, 1\}$.
836 We need to encode the x 's into W somehow: We aim at the following embedding:

$$W = \begin{pmatrix} x_1 & \bar{x}_1 & \dots & x_n & \bar{x}_n & y_0 & \dots & y_k \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1 & \bar{x}_1 & \dots & x_n & \bar{x}_n & y_0 & \dots & y_k \end{pmatrix}$$

837 The y are $k+1 := \max\{1, m-n+2\}$ extra dummy variables to make the matrix a square $d \times d$ matrix
838 with $d := \max\{m+n+2, 2n+1\}$.

839 Choosing a cyclic permutation matrix $\Pi = [234\dots d1]$ ensures that all rows of W are indeed the same
840 via $\Pi \cdot W = W$. The standard way of achieving $x_j, y_j \in \{0, 1\}$ is via $x_j^2 = x_j$ and $y_j^2 = y_j$. This can be
841 achieved via $(\text{Id} \odot W)^2 = \text{Id} \odot W$, were $\text{Id}_{ss'} = \delta_{ss'}$ is the identity matrix.

842 We use $[B \odot W]_{s+} = 1$ to ensure $\bar{x}_j = 1 - x_j$, $y_0 = 1$, and $y_1 = \dots = y_k = 0$ and $\ell_{i1} + \ell_{i2} + \ell_{i3} = 1$ by
843 setting $B_{s, 2s-1} = B_{s, 2s} = 1$ for $s \in \{1:n\}$, and $B_{i+n, 2j-1} = 1$ if $\ell_{ia} = x_j$ and $B_{i+n, 2j} = 1$ if $\ell_{ia} = \neg x_j$
844 for $i \in \{1:m\}$ and $a \in \{1, 2, 3\}$, and $B_{d-1, 2n+1} = \dots = B_{d-1, 2n+m} = 1$, and $B_{d, 2n+1} = 1$, and $B_{ss'} = 0$
845 for all other ss' . This also ensures that all rows of W sum to $n+1$, hence $W \cdot W = (n+1)W$, so
846 $x_j \in \{0, 1\}$ can be achieved via $C = \text{Id}$ and $A = \frac{1}{n+1} \text{Id}$ in $A \odot (W \cdot W) = (C \odot W)(C \odot W)$.

847 The construction implies that the 3CNF(n, m, g) formula is satisfiable iff (33) has a solution in W
848 with the A, B, C, Π as constructed above. This shows NP-hardness of deciding whether (33) has a
849 solution. A solution can trivially be verified (in the rationals or to ε -precision over the reals) in time
850 $O(d^3)$, hence the problem is in NP, hence NP-complete.

851 L Compact Representation of EqIM(2+)

852 If only B^{a+} (EqIM(2+)) is given, we can sum (32) over a' . If we further assume $a=2$ and define
853 $B = B^0$ and $A = B^{0+}$ and $W = W^+$ and exploit $B^+ = B^{++} = 1$, this reduces to the elegant quadratic
854 matrix equation

$$A \odot (W \cdot W) = (B \odot W) \cdot W \quad (34)$$

855 with constraints as in (32), or even simpler $W_{s+} = 1$ if π is unknown. This is the most pure formulation
856 of the problem we are trying but are unable to solve we could come up with. For A and B defined via
857 M , we know that (34) has a solution (namely $W = M^+$).

858 We neither know whether there exists an efficient algorithm to find *some* solution (34), nor to find *the*
859 solution in case it is unique, nor to decide whether there exist solutions in case A and B are chosen
860 arbitrarily.

861 The condition $W_{s+} = 1$ can be relaxed to $W_{s+} > 0$. If $W_{ss'}$ is a solution of (34), then also $v_s^{-1} W_{ss'} v_{s'}$
862 for any $v_s > 0$ (most easily checked via (11)). Every non-negative matrix has a real non-negative
863 Eigenvector v , and $W_{s+} > 0$ implies $v_s > 0$ and Eigenvalue $\lambda > 0$, hence for $W_{ss'}^{\text{norm}} := (\lambda v_s)^{-1} W_{ss'} v_{s'}$,
864 we have $W_{s+}^{\text{norm}} = 1$.

865 $B^a \geq 0$ and $B^+ = 1$ iff $B \in [0; 1]$ (and $B^1 = 1 - B$). $B^{a+} \geq 0$ and $B^{++} = 1$ iff $A \in [0; 1]$ (and
866 $B^{1+} = 1 - A$). But we can scale back any A and B by the same $0 < \lambda < 1$ to satisfy these without
867 changing (34), i.e. these extra conditions (A and B bounded by 1) do not make the problem any
868 simpler.

869 M Open Problem

870 We present the most important open problem(s) in their simplest instantiation and most elegant form,
871 fully self-contained here: Consider matrices $A, B, W \in [0; 1]^{d \times d}$ with $d \in \mathbb{N}$, tied by the quadratic
872 matrix equation

$$A \odot (W \cdot W) = (B \odot W) \cdot W \quad \text{and} \quad W_{s+} = 1 \forall s \quad (35)$$

873 where \odot is element-wise (Hadamard) multiplication and \cdot is standard matrix multiplication. The open
874 problems are as follows: Given A and B , are there efficient algorithms which

- 875 (a) decide whether there exists a W satisfying (35)?
- 876 (b) decide whether the solution is unique, assuming (35) has a solution?
- 877 (c) compute a solution, assuming (35) has a solution?
- 878 (d) compute *the* solution, assuming (35) has a unique solution?

879 Computing a real number means, given any $\varepsilon > 0$, computing an ε -approximation. Efficient means
 880 running time is polynomial in d , ideally with a degree independent of $1/\varepsilon$. General systems of
 881 quadratic equations are known to be NP-hard, but we do not know the complexity of this particular
 882 matrix sub-class.

883 The upper bounds $A, B, W \leq 1$ can always be satisfied by scaling, hence are irrelevant. $W_{s+} = 1$ can be
 884 relaxed to $W_{s+} > 0$ except in the uniqueness questions. If helpful: One may assume A, B, W strictly
 885 positive. Also, any finite (d -independent) number of equations of the form $A' \odot (W \cdot W) = (B' \odot W) \cdot$
 886 W with other *general* matrices $A', B' \in [0; 1]^{d \times d}$ may be added, which further constrain the solution
 887 space.

888 N Further Experiments

889 Here we provide further experiments supporting and illustrating the theory. In Appendix O we show
 890 how we numerically dealt with $B = 0/0 = \perp$. Appendix P derives the formulas for the plotted solution
 891 dimensions.

892 **Experiments illustrating robustness to noise.** The propositions and results in the main text assume
 893 that we know the one and two step inverse models ($B1 := B^a$, $B2 := B^{a+}$) exactly, but in practice
 894 these distributions must be estimated from data. Here we investigate the extent to which our algorithm
 895 is robust to noise arising from learning.

896 Rather than committing to a specific learning algorithm, we instead directly inject noise into the true
 897 inverse distributions. This is done by adding $\varepsilon \times 10^c$ to the true distribution and renormalizing B ,
 898 where ε is drawn from the unit uniform distribution: $\varepsilon \sim \mathcal{U}[0, 1]$

899 Figure 4a shows that noise doesn't substantially degrade performance across several orders of
 900 magnitude (c varied -7 to 0). Additionally, the effect of this noise is substantially diminished as the
 901 horizon of the inverse model is increased (from $B1 := B^a$ to $B3 := B^{a++}$). While the is perhaps not
 902 surprising, as the entropy of such inverse distributions increases monotonically with the horizon, it
 903 still shows that noise is not compounding in a way that renders long-horizon predictions meaningless.
 904 Figure 2 buttresses this interpretation by showing that the recovered B^{a++} is qualitatively similar to
 905 the ground truth even with substantial noise.

906 **Experiments on the Tensor-product special case.** As detailed in Section 3, if M factors into two
 907 processes $\dot{M}^a \otimes \dot{M}$, where \dot{M} is action-independent, then only the complexity of the action-dependent
 908 process \dot{M}^a matters for all of our questions.

909 This particular special case is important because of its frequency in applied work. Many environments
 910 have most of their complexity in sub-spaces that the agent has no control over. This is illustrated by
 911 Figure 3, reproduced from [LFLDP21], wherein naturalistic videos are superimposed on relatively
 912 simple continuous control environments. Clearly, the background dynamics can be arbitrarily complex
 913 without impacting the underlying control problem.

914 We can construct small environments of this form via a simple procedure. We construct \dot{M} with \dot{d}
 915 states and k actions by sampling each element of the appropriately sized matrices from $\mathcal{U}[0, 1]$ and
 916 then normalizing. \dot{M} has two states that transition uniformly regardless of the action.

917 The linear algorithm of Section 4 can (implicitly) output all W and $B2$ consistent with $B1$, and the
 918 formulas derived in Appendix P allow to (explicitly) calculate the dimensions of the solution spaces.

919 In the experiments shown in Figure 4b, $k = 5$ as in the main text, and $d = 2\dot{d}$ is varied from 16 to 32.
 920 The results show that the space of forward dynamics W is always significantly larger than the space
 921 of the 2-step inverse models ($B2$). This confirms that inverse models can be significantly simpler
 922 than forward models.

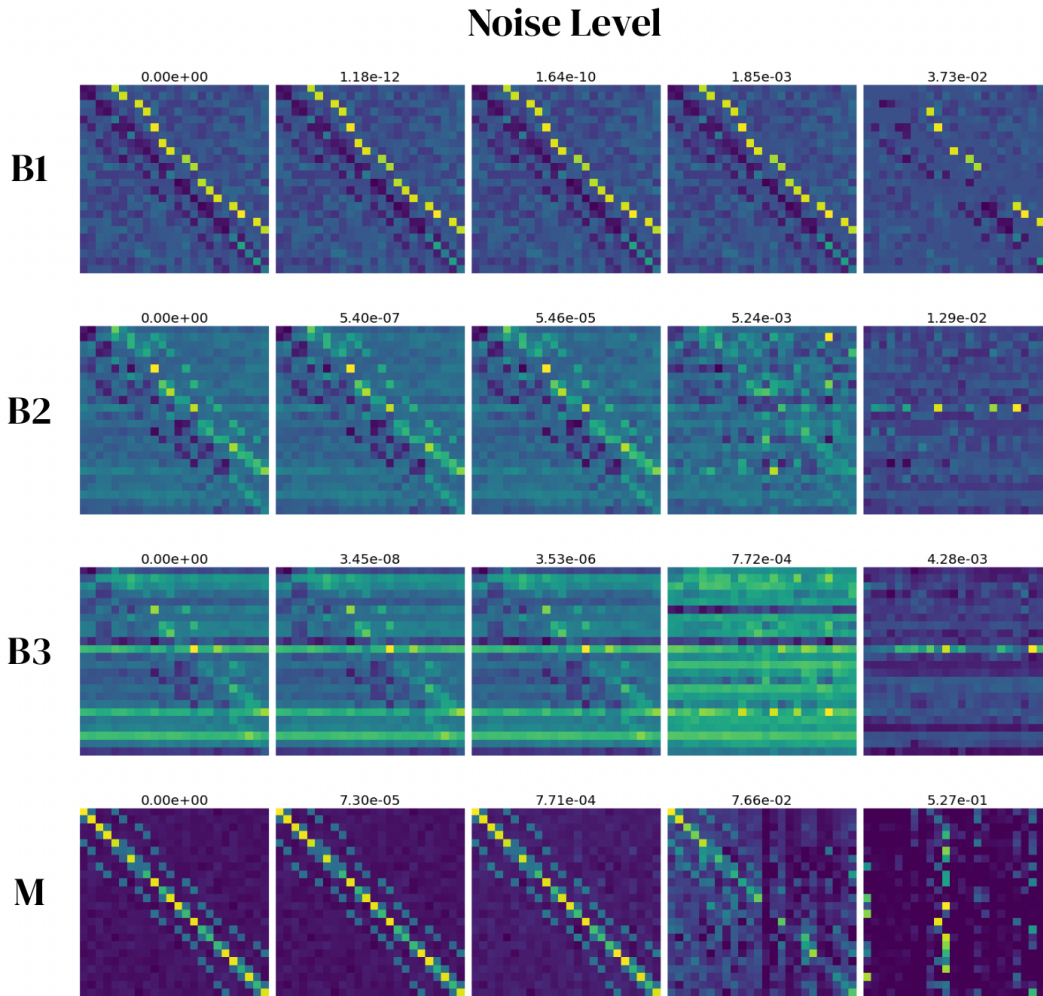
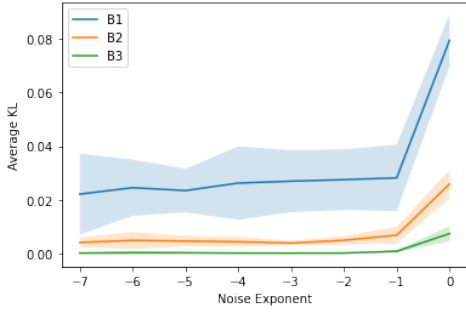


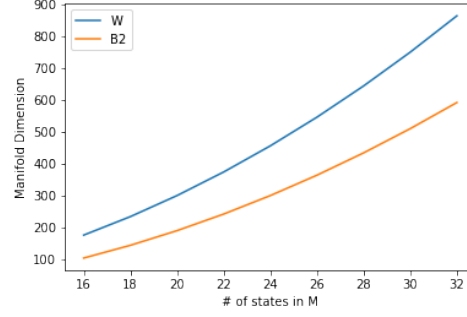
Figure 2: Reconstructing inverse and forward models from inverse models with noise injected. Noise increases exponentially across columns $[0, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}]$. The subplot titles show the average KL divergence of the recovered distribution from the ground truth.



Figure 3: Reproduced from [LFLDP21], this 'half-cheetah' environment has been augmented with videos of complex scenes. This highlights how non-controllable aspects of the environment can be made more complex without changing the underlying control problem. The fact that such environments are of interest motivates our focus on the Tensor-product special case.



(a) Effect of Noise



(b) Tensor-product Solution Space

Figure 4: **(a) Noise-induced reconstruction error:** In practice W must be inferred from learned estimates of $B1$ and $B2$. We investigate the effect of the resulting error on the inverse models ($B1, B2, B3$) recovered from the inferred W in terms of their proximity to the ground truth distributions. At each noise level the algorithm was run on 10 randomly generated grids, with the shaded region representing $\pm 2\sigma$. **(b) Solution dimensions of W and $B2$ given $B1$:** When the solution to an inverse model ($B2$) given only $B1$ is not unique, we can characterize the solution space in terms of its manifold dimension. By comparing this to the dimension of that of the inferred forward model (W), we can see that our algorithm has narrowed down the space of inverse models significantly more. If also $B2$ is given, the solution dimension of W reduces from d_W (blue curve) to $d_W - d_B$ (blue-orange curve).

923 O How to Deal with 0/0

924 If for some pair of states (s, s') , no action a of positive π -probability leads from state s to s' , i.e. if
 925 $M_{ss'}^+ = 0$, then $B_{ss'}^+$ and $B_{ss'}^a, \forall a$ are $0/0 = \perp = \text{undefined}$. To also handle $B_{ss'} = \perp$, we need to adapt
 926 the linear algorithm in Section 4. We provide 2 different ways of doing so, with a couple of variations,
 927 all leading to the same correct result.

928 We have to restrict the sum in $\sum_{s'} B_{ss'}^a J_{ss'} = \pi(a|s)$ to those s' for which $B_{ss'}^a$ is defined. We then
 929 solve for $J_{ss'}$, again for s' for which $B_{ss'}^a$ is defined, and set $J_{ss'} = 0$ for those s' for which $B_{ss'}^a = \perp$.
 930 Technically this can be achieved by removing the s' columns from matrix B_s^a and J_s for which
 931 $B_{ss'}^a = \perp$, solve the reduced linear equation system, and finally reinsert $J_{ss'} = 0$ for the removed s' .
 932 Simpler is to replace $B_{ss'}^a = \perp$ by $B_{ss'}^a = 0$, solve the equation for J , and then set $J_{ss'} = 0$ for the s'
 933 for which the original $B_{ss'}^a$ was \perp . Some solvers automatically result in $J_{ss'} = 0$, since this is the
 934 minimum norm solution, but it is better not to rely on this. Instead of setting $J_{ss'} = 0$ after solving
 935 the linear system, one could also augment B_s^a with extra rows that enforce $J_{ss'} = 0$.

936 Alternatively, we could replace $B_{ss'} = \perp$ by a random vector which sums to 1, e.g. $B_{ss'}^a = r_a / r_+$,
 937 where $r_a = -\log u_a$ with $u_a \sim \text{Uniform}[0;1]$. Provided that the solution is unique, this also leads to the
 938 correct solution (almost surely), and in this way $J_{ss'} = 0$ automatically. If the solution is not unique,
 939 W will still satisfy $B^a = W^a \circledast W^+$ when for $B_{ss'}^a \neq \perp$, but $W_{ss'}$ may not be 0.

940 The adaptation of the Linear Relaxation Algorithm in Section 5 follows the same pattern: $A_{ss^i s^j} = \perp$
 941 in (12), whenever one of the three involved B 's is undefined. For such $ss^i s^j$, we need to ensure that
 942 $\hat{U}_{ss^i s^j} = 0$, which can be done with any of the variations described above. Once we have $\hat{U}_{ss^i s^j}$, we
 943 set $C_{ss^i s^j}^a = 0$ if $B_{ss^i s^j}^a = \perp$. No further intervention is needed, since $\hat{U}_{ss^i s^j} = 0$ already.

944 P Solution Dimensions of W and $B^{aa'}$.

945 In Section 4 we presented an algorithm for inferring W and $B^{aa'}$ from B^a . Even if M cannot
 946 uniquely be reconstructed \neg (i), $B^{aa'}$ may still be unique (iii). More generally, the solutions J and
 947 W^a form linear spaces of dimension $d_J = d_W \leq d(d-1)$ ($d_J \geq d_W$ since W is linear function of
 948 J . $d_J \leq d_W$, since $W^+ = J$). $B^{aa'}$ is a (non-linear, polynomial) variety of dimension $d_B \leq d_W$ at
 949 regular points (since it is a smooth function of W).

950 **Parameterizing the solutions for J and W and B .** We can determine the solution dimensions
951 d_J , d_W , and d_B as follows: Let $Y_{ss'}$ be a solution of $[B^a \odot Y]_{s+} = 0$. If $\hat{J}_{ss'}$ is a solution of
952 $[B^a \odot J]_{s+} = \pi(a|s)$, then so is $J := \hat{J} + Y$, hence $W^a := \hat{W}^a + X^a$ is a solution of $B^a = W^a \odot W^+$
953 and $W_{s+}^a = \pi(a|s)$, where $\hat{W}^a := B^a \odot \hat{J}$ and $X^a := B^a \odot Y$.

954 If we plug in $W^a \equiv \hat{W}^a + X^a$ into $B^{aa'}$, we get the variety of $B^{aa'}$ parameterized in terms X^a . If
955 we expand this non-linear expression up to linear order in X^a , we get after some algebra

$$B^{aa'} = [\hat{W}^a \hat{W}^{a'} + \hat{W}^a X^{a'} + X^a \hat{W}^{a'} - (\hat{W}^a \hat{W}^{a'}) \odot (\hat{W}^+)^2 \odot (\hat{W}^+ X^+ + X^+ \hat{W}^+)] \odot (\hat{W}^+)^2 + O(X^2) \quad (36)$$

956 The linear part forms a tangent direction on the variety at $\hat{B}^{aa'}$.

957 **Determining the solution dimensions for J and W and B .** Now, for each s , let $Y_{ss'}^r$ for
958 $r \in \{1 : d_{J_s}\}$ span all solutions of $[B^a \odot Y]_{s+} = 0$, which can easily be determined by SVD: d_{J_s} is
959 the number zero singular values of matrix B_{s+}^a , and Y_s^r the corresponding singular vectors. Then,
960 $J_{ss'} = \hat{J}_{ss'} + \sum_r Y_{ss'}^r z_{sr}$ for any $z \in \mathbb{R}^{d_J}$ with $d_J = \sum_s d_{J_s}$ is a solution of $[B^a \odot J]_{s+} = \pi(a|s)$.

961 Similarly, $W_{ss'}^a := \hat{W}_{ss'}^a + \sum_r X_{ss'}^{ar} z_{sr}$ with $X^{ar} := B^a \odot Y^r$ span all solutions consistent with B^a and
962 π . The solution dimension is $d_W = \sum_s d_{W_s}$, where for each s , d_{W_s} is the rank of X_s^a : if interpreted
963 as a $kd \times d_{J_s}$ matrix in $as' \times r$. d_{W_s} may be smaller than d_{J_s} , since unlike Y_s^r , X_s^a may not be full
964 rank.

965 If we plug $X_{ss'}^a = \sum_r X_{ss'}^{ar} z_{sr}$ into (36), after some index manipulation we get

$$B^{aa'} = \hat{B}^{aa'} + \sum_{t=1}^d \sum_{r=1}^{d_{J_t}} C^{aa'rt} z_{tr} \odot (\hat{W}^+)^2 \odot (\hat{W}^+)^2 + O(z^2) \quad \text{with} \quad (37)$$

$$C_{ss''}^{aa'rt} := (\hat{W}_{st}^a X_{ts''}^{a'r} + [X^{ar} \hat{W}^{a'}]_{ss''} \delta_{ts}) [(\hat{W}^+)^2]_{ss''} - [\hat{W}^a \hat{W}^{a'}]_{ss''} (\hat{W}_{st}^+ X_{ts''}^{+r} + [X^{+r} \hat{W}^+]_{ss''} \delta_{ts}) \quad (38)$$

966 $B^{aa'}(z)$ is a local parametrization of B , and if we drop the $O(z^2)$, it parameterize its tangential
967 hyperplane at $\hat{B}^{aa'}$. Its dimension d_B is the rank of C interpreted as a $k^2 d^2 \times d_J$ matrix in $aa' ss'' \times rt$.
968 Again, d_B may be smaller than d_W , since C may not be full rank. for $r \in \{1 : d_J\}$ spans the
969 tangential space of rescaled variety $B^{aa'}$ at $\hat{B}^{aa'}$. Again, they may not be linearly independent, If
970 $[(W^+)^2]_{ss''} = 0$, then $B_{ss''}^{aa'} = \perp \forall aa'$, hence all such ss'' should be ignored in $C_{ss''}^{aa'rt}$, but since the
971 corresponding rows in C are 0, they don't contribute to the rank anyway.

972 **Sampling estimate of d_B .** A simpler, but less elegant, and more fragile method to estimate d_B is
973 as follows: Fix one solution \hat{J} . Add random noise in direction of the null-space spanned by Y^r so
974 that it stays a solution, i.e. compute $J = \hat{J} + \sum_r Y^r z_r$ for random z , and from this, W and $B^{aa'}$ for
975 many such random J . The resulting point cloud spans covers the solution variety $B^{aa'}$. Various tools
976 could be used to analyze this point cloud, e.g. determine its dimension. If z is chosen small, the point
977 cloud concentrates around $\hat{B}^{aa'}$ and forms a near-linear space, whose dimension d_B can easily be
978 determined by PCA.

979 **Higher-order B and higher i .** In the same way we can derive the solution dimensions $d_{B^{\dots}}$ for
980 higher-order B^{\dots} . Also, even though we don't have (yet) an efficient algorithm for solving EqIM(i)
981 for $i > 1$ if the solution is not unique, we still can determine the dimension of the solutions (at a
982 particular point \hat{W}). Algorithmically already covered is the case of W satisfying EqIM(1) \wedge EqIM(2),
983 whose solution dimension turns out to be $d_W - d_B$. The general procedure is to plug $W = \hat{W} + X$
984 into and linearly expand EqIM(i) for i we to hold. Together they form a system of linear equations
985 whose solution dimension can be determined by SVD as above.

986 Q Counter-Examples in Related Work

987 In Section 3 we presented a counter-example to questions (i,iii,v). Question (i) (i.e. Can M be
988 inferred from $B^a := M^a \odot M^+$?) has been implicitly addressed in previous work. In [EMK⁺22,
989 App.A.3] the authors present a counter-example to the claim that a state representation constructed

990 via an inverse model (i.e. two states have the same representation iff they yield the same inverse
991 distribution for all of their possible successor states) is sufficient for representing a set of policies that
992 differentially visit all states¹. This fails whenever two states are aliased by the inverse model.

993 Note that this failure of state representation learning implies a negative answer to our question (i),
994 as W would differ from M on these aliased states. Unlike our counter-example, theirs involves
995 deterministic forward dynamics, and therefor buttresses our claims by showing that M cannot always
996 be inferred even in this simpler case. Similar to our counter-example in Section 3, [MHKL20]
997 proposes a stochastic counter-example to inverse modeling for state representation learning.

998 In general, the transferability of these counter-examples suggests a strong relationship between
999 literature on using single-step inverse models for state representation learning and using them
1000 for inferring the forward model. It is an interesting open question whether or not algorithms
1001 for representation learning on the basis on multi-step inverse models (like those put forward in
1002 [EMK⁺22]) might be used to shed light on the questions put forward here and vice versa.

1003 **R Relevance for Planning**

1004 In Section 1, various streams of applied work were highlighted; here we focus on spelling out the
1005 overarching impact that compositional inverse models (an affirmative answer to question (iv)) would
1006 have for planning problem.

1007 Many forms of planning involve the evaluation of candidate i -step action sequences (e.g. model
1008 predictive path integral control [WDG⁺16]). Ideally, all possible action sequences would be evaluated,
1009 but as the space of i -step action sequences grows exponentially in i , this is often intractable.

1010 Access to the i -step inverse distribution $p(a...a^i | s...s^{i+1})$ allows determining the subset of action
1011 sequences that likely reach state s^{i+1} post-execution (e.g. those whose probability is above some
1012 threshold). It is often the case that only action sequences that are distinguishable in this way are of
1013 interest (e.g. goal-reach tasks), thus access to an inverse model of the appropriate horizon allows for
1014 filtering candidates. This filtering method is a particularly appealing approach when the cost/reward
1015 function is initially unknown and frequently changes.

1016 While this idea has already seen scalable implementations [MJR15], these rely on short, fixed horizons
1017 since they directly learn all inverse models of step-size up to the horizon, which is data-inefficient for
1018 large horizons. If inverse models could be composed, then longer, variable horizons could be used
1019 while only learning a short horizon inverse model by inferring the longer horizon models as needed.
1020 Our work shows that this is possible, but with exceptions and a more practical composition algorithm
1021 being outstanding.

¹Technically, as per their Definition 2, this ‘policy cover’ need only account for all ‘endogenous’ states. But omit the ‘exogenous’ states from their counter-example and it can be seen to address our question (i).