

Verified Algorithms for Solving Markov Decision Processes

Maximilian Schöffeler and Mohammad Abdulaziz

December 28, 2021

Abstract

We present a formalization of algorithms for solving Markov Decision Processes (MDPs) with formal guarantees on the optimality of their solutions. In particular we build on our analysis of the Bellman operator for discounted infinite horizon MDPs. From the iterator rule on the Bellman operator we directly derive executable value iteration and policy iteration algorithms to iteratively solve finite MDPs. We also prove correct optimized versions of value iteration that use matrix splittings to improve the convergence rate. In particular, we formally verify Gauss-Seidel value iteration and modified policy iteration. The algorithms are evaluated on two standard examples from the literature, namely, inventory management and gridworld. Our formalization covers most of chapter 6 in Puterman’s book [1].

Contents

1	Value Iteration	2
2	Policy Iteration	6
3	Modified Policy Iteration	9
3.1	The Advantage Function B	10
3.2	Optimization of the Value Function over Multiple Steps	10
3.3	Expressing a Single Step of Modified Policy Iteration	11
3.4	Computing the Bellman Operator over Multiple Steps	12
3.5	The Modified Policy Iteration Algorithm	12
3.6	Convergence Proof	13
3.7	ϵ -Optimality	14
3.8	Unbounded MPI	14
3.9	Initial Value Estimate $v_0\text{-mpi}$	16
3.10	An Instance of Modified Policy Iteration with a Valid Conservative Initial Value Estimate	16

4	Matrices	17
4.1	Nonnegative Matrices	17
4.2	Matrix Powers	17
4.3	Triangular Matrices	18
4.4	Inverses	19
5	Bounded Linear Functions and Matrices	19
6	Value Iteration using Splitting Methods	22
6.1	Regular Splittings for Matrices and Bounded Linear Functions	22
6.2	Splitting Methods for MDPs	23
6.3	Discount Factor <i>QR-disc</i>	24
6.4	Bellman-Operator	24
6.5	Gauss Seidel Splitting	26
6.5.1	Definition of Upper and Lower Triangular Matrices	26
6.5.2	Gauss Seidel is a Regular Splitting	27
7	Code Generation for MDP Algorithms	37
7.1	Least Argmax	37
7.2	Functions as Vectors	39
7.3	Bounded Functions as Vectors	39
7.4	IArrays with Lengths in the Type	40
7.5	Value Iteration	41
7.6	Policy Iteration	42
7.7	Gauss-Seidel Iteration	43
7.8	Modified Policy Iteration	44
7.9	Auxiliary Equations	45
8	Code Generation for Concrete Finite MDPs	46
9	Inventory Management Example	48
10	Gridworld Example	51

```

theory Value-Iteration
imports MDP-Rewards.MDP-reward
begin

context MDP-att-L
begin

```

1 Value Iteration

In the previous sections we derived that repeated application of \mathcal{L}_b to any bounded function from states to the reals converges to the optimal value of the MDP ν_b-opt .

We can turn this procedure into an algorithm that computes not only an approximation of $\nu_b\text{-opt}$ but also a policy that is arbitrarily close to optimal.

Most of the proofs rely on the assumption that the supremum in \mathcal{L}_b can always be attained.

The following lemma shows that the relation we use to prove termination of the value iteration algorithm decreases in each step. In essence, the distance of the estimate to the optimal value decreases by a factor of at least l per iteration.

lemma *vi-rel-dec*:

assumes $l \neq 0 \ \mathcal{L}_b \ v \neq \nu_b\text{-opt}$

shows $\lceil \log (1 / l) (\text{dist } (\mathcal{L}_b \ v) \ \nu_b\text{-opt}) - c \rceil < \lceil \log (1 / l) (\text{dist } v \ \nu_b\text{-opt}) - c \rceil$

<proof>

lemma *dist- \mathcal{L}_b -lt-dist-opt*: $\text{dist } v \ (\mathcal{L}_b \ v) \leq 2 * \text{dist } v \ \nu_b\text{-opt}$

<proof>

abbreviation *term-measure* $\equiv (\lambda(\text{eps}, v).$

if $v = \nu_b\text{-opt} \vee l = 0$

then 0

else $\text{nat } (\text{ceiling } (\log (1/l) (\text{dist } v \ \nu_b\text{-opt}) - \log (1/l) (\text{eps} * (1-l) / (8 * l))))$

function *value-iteration* $:: \text{real} \Rightarrow ('s \Rightarrow_b \text{real}) \Rightarrow ('s \Rightarrow_b \text{real})$ **where**

value-iteration eps v =

(if $2 * l * \text{dist } v \ (\mathcal{L}_b \ v) < \text{eps} * (1-l) \vee \text{eps} \leq 0$ *then* $\mathcal{L}_b \ v$ *else* *value-iteration eps* $(\mathcal{L}_b \ v)$)

<proof>

termination

<proof>

The distance between an estimate for the value and the optimal value can be bounded with respect to the distance between the estimate and the result of applying it to \mathcal{L}_b

lemma *contraction- \mathcal{L} -dist*: $(1 - l) * \text{dist } v \ \nu_b\text{-opt} \leq \text{dist } v \ (\mathcal{L}_b \ v)$

<proof>

lemma *dist- \mathcal{L}_b -opt-eps*:

assumes $\text{eps} > 0 \ 2 * l * \text{dist } v \ (\mathcal{L}_b \ v) < \text{eps} * (1-l)$

shows $\text{dist } (\mathcal{L}_b \ v) \ \nu_b\text{-opt} < \text{eps} / 2$

<proof>

The estimates above allow to give a bound on the error of *value-iteration*.

declare *value-iteration.simps*[*simp del*]

lemma *value-iteration-error*:

assumes $eps > 0$

shows $dist (value\text{-}iteration\ eps\ v)\ \nu_b\text{-}opt < eps / 2$

<proof>

After the value iteration terminates, one can easily obtain a stationary deterministic epsilon-optimal policy.

Such a policy does not exist in general, attainment of the supremum in \mathcal{L}_b is required.

definition *find-policy* ($v :: 's \Rightarrow_b real$) $s = arg\text{-}max\text{-}on (\lambda a. L_a\ a\ v\ s)$
($A\ s$)

definition *vi-policy* $eps\ v = find\text{-}policy (value\text{-}iteration\ eps\ v)$

We formalize the attainment of the supremum using a predicate *has-arg-max*.

abbreviation $vi\ u\ n \equiv (\mathcal{L}_b \ \widetilde{\sim}^n)\ u$

lemma *\mathcal{L}_b -iter-mono*:

assumes $u \leq v$ **shows** $vi\ u\ n \leq vi\ v\ n$

<proof>

lemma

assumes $vi\ v\ (Suc\ n) \leq vi\ v\ n$

shows $vi\ v\ (Suc\ n + m) \leq vi\ v\ (n + m)$

<proof>

lemma

assumes $vi\ v\ n \leq vi\ v\ (Suc\ n)$

shows $vi\ v\ (n + m) \leq vi\ v\ (Suc\ n + m)$

<proof>

lemma $vi\ v \longrightarrow \nu_b\text{-}opt$

<proof>

lemma $(\lambda n. dist (vi\ v\ (Suc\ n)) (vi\ v\ n)) \longrightarrow 0$

<proof>

end

context *MDP-att- \mathcal{L}*

begin

The error of the resulting policy is bounded by the distance from its value to the value computed by the value iteration plus the error in the value iteration itself. We show that both are less than $\text{eps} / (2::'b)$ when the algorithm terminates.

lemma *find-policy-error-bound*:

assumes $\text{eps} > 0$ $2 * l * \text{dist } v (\mathcal{L}_b v) < \text{eps} * (1-l)$
shows $\text{dist } (\nu_b (\text{mk-stationary-det } (\text{find-policy } (\mathcal{L}_b v)))) \nu_b\text{-opt} < \text{eps}$
 $\langle \text{proof} \rangle$

lemma *vi-policy-opt*:

assumes $0 < \text{eps}$
shows $\text{dist } (\nu_b (\text{mk-stationary-det } (\text{vi-policy } \text{eps } v))) \nu_b\text{-opt} < \text{eps}$
 $\langle \text{proof} \rangle$

lemma *lemma-6-3-1-d*:

assumes $\text{eps} > 0$
assumes $2 * l * \text{dist } (vi v (\text{Suc } n)) (vi v n) < \text{eps} * (1-l)$
shows $\text{dist } (vi v (\text{Suc } n)) \nu_b\text{-opt} < \text{eps} / 2$
 $\langle \text{proof} \rangle$

end

context *MDP-act* **begin**

definition *find-policy'* ($v :: 's \Rightarrow_b \text{real}$) $s = \text{arb-act } (\text{opt-acts } v s)$

definition *vi-policy'* $\text{eps } v = \text{find-policy}' (\text{value-iteration } \text{eps } v)$

lemma *find-policy'-error-bound*:

assumes $\text{eps} > 0$ $2 * l * \text{dist } v (\mathcal{L}_b v) < \text{eps} * (1-l)$
shows $\text{dist } (\nu_b (\text{mk-stationary-det } (\text{find-policy}' (\mathcal{L}_b v)))) \nu_b\text{-opt} < \text{eps}$
 $\langle \text{proof} \rangle$

lemma *vi-policy'-opt*:

assumes $\text{eps} > 0$ $l > 0$
shows $\text{dist } (\nu_b (\text{mk-stationary-det } (\text{vi-policy}' \text{eps } v))) \nu_b\text{-opt} < \text{eps}$
 $\langle \text{proof} \rangle$

end

end

theory *Policy-Iteration*

imports *MDP-Rewards.MDP-reward*

begin

2 Policy Iteration

The Policy Iteration algorithms provides another way to find optimal policies under the expected total reward criterion. It differs from Value Iteration in that it continuously improves an initial guess for an optimal decision rule. Its execution can be subdivided into two alternating steps: policy evaluation and policy improvement.

Policy evaluation means the calculation of the value of the current decision rule.

During the improvement phase, we choose the decision rule with the maximum value for L , while we prefer to keep the old action selection in case of ties.

context *MDP-att- \mathcal{L}* **begin**

definition *policy-eval* $d = \nu_b$ (*mk-stationary-det* d)
end

context *MDP-act*

begin

definition *policy-improvement* $d \ v \ s =$ (
 if is-arg-max ($\lambda a. L_a \ a$ (*apply-bfun* v) s) ($\lambda a. a \in A \ s$) ($d \ s$)
 then $d \ s$
 else arb-act (*opt-acts* $v \ s$)

definition *policy-step* $d = \text{policy-improvement } d \ (\text{policy-eval } d)$

function *policy-iteration* $:: ('s \Rightarrow 'a) \Rightarrow ('s \Rightarrow 'a)$ **where**
 policy-iteration $d =$ (
 let $d' = \text{policy-step } d$ *in*
 if $d = d' \vee \neg \text{is-dec-det } d$ *then* d *else* *policy-iteration* d')
 $\langle \text{proof} \rangle$

The policy iteration algorithm as stated above does require that the supremum in \mathcal{L}_b is always attained.

Each policy improvement returns a valid decision rule.

lemma *is-dec-det-pi*: *is-dec-det* (*policy-improvement* $d \ v$)
 $\langle \text{proof} \rangle$

lemma *policy-improvement-is-dec-det*: $d \in D_D \implies \text{policy-improvement } d \ v \in D_D$
 $\langle \text{proof} \rangle$

lemma *policy-improvement-improving*:
 assumes $d \in D_D$

shows ν -improving v ($mk\text{-}dec\text{-}det$ ($policy\text{-}improvement$ d v))
 $\langle proof \rangle$

lemma *eval-policy-step-L*:

assumes $is\text{-}dec\text{-}det$ d

shows L ($mk\text{-}dec\text{-}det$ ($policy\text{-}step$ d)) ($policy\text{-}eval$ d) = \mathcal{L}_b ($policy\text{-}eval$ d)
 $\langle proof \rangle$

The sequence of policies generated by policy iteration has monotonically increasing discounted reward.

lemma *policy-eval-mon*:

assumes $is\text{-}dec\text{-}det$ d

shows $policy\text{-}eval$ $d \leq policy\text{-}eval$ ($policy\text{-}step$ d)

$\langle proof \rangle$

If policy iteration terminates, i.e. $d = policy\text{-}step$ d , then it does so with optimal value.

lemma *policy-step-eq-imp-opt*:

assumes $is\text{-}dec\text{-}det$ d $d = policy\text{-}step$ d

shows ν_b ($mk\text{-}stationary$ ($mk\text{-}dec\text{-}det$ d)) = $\nu_b\text{-}opt$

$\langle proof \rangle$

end

We prove termination of policy iteration only if both the state and action sets are finite.

locale *MDP-PI-finite* = *MDP-act* A K r l *arb-act*

for

A **and**

$K :: 's :: countable \times 'a :: countable \Rightarrow 's$ *pmf* **and** r l *arb-act* +

assumes *fin-states*: $finite$ ($UNIV :: 's$ *set*) **and** *fin-actions*: $\bigwedge s.$ $finite$ (A s)

begin

If the state and action sets are both finite, then so is the set of deterministic decision rules D_D

lemma *finite- D_D [simp]*: $finite$ D_D

$\langle proof \rangle$

lemma *finite-rel*: $finite$ $\{(u, v). is\text{-}dec\text{-}det$ $u \wedge is\text{-}dec\text{-}det$ $v \wedge \nu_b$ ($mk\text{-}stationary\text{-}det$ u) >

ν_b ($mk\text{-}stationary\text{-}det$ v)}

$\langle proof \rangle$

This auxiliary lemma shows that policy iteration terminates if no improvement to the value of the policy could be made, as then the policy remains unchanged.

lemma *eval-eq-imp-policy-eq*:

assumes *policy-eval* $d = \text{policy-eval } (\text{policy-step } d)$ *is-dec-det* d

shows $d = \text{policy-step } d$

<proof>

We are now ready to prove termination in the context of finite state-action spaces. Intuitively, the algorithm terminates as there are only finitely many decision rules, and in each recursive call the value of the decision rule increases.

termination *policy-iteration*

<proof>

The termination proof gives us access to the induction rule/simplification lemmas associated with the *policy-iteration* definition. Thus we can prove that the algorithm finds an optimal policy.

lemma *is-dec-det-pi'*: $d \in D_D \implies \text{is-dec-det } (\text{policy-iteration } d)$

<proof>

lemma *pi-pi[simp]*: $d \in D_D \implies \text{policy-step } (\text{policy-iteration } d) = \text{policy-iteration } d$

<proof>

lemma *policy-iteration-correct*:

$d \in D_D \implies \nu_b (\text{mk-stationary-det } (\text{policy-iteration } d)) = \nu_b\text{-opt}$

<proof>

end

context *MDP-finite-type* **begin**

The following proofs concern code generation, i.e. how to represent \mathcal{P}_1 as a matrix.

sublocale *MDP-att-L*

<proof>

definition *fun-to-matrix* $f = \text{matrix } (\lambda v. (\chi j. f (\text{vec-nth } v) j))$

definition *Ek-mat* $d = \text{fun-to-matrix } (\lambda v. ((\mathcal{P}_1 d) (\text{Bfun } v)))$

definition *nu-inv-mat* $d = \text{fun-to-matrix } ((\lambda v. ((\text{id-blifun} - l *_R \mathcal{P}_1 d) (\text{Bfun } v))))$

definition *nu-mat* $d = \text{fun-to-matrix } (\lambda v. ((\sum i. (l *_R \mathcal{P}_1 d) \overset{\sim}{\sim} i) (\text{Bfun } v)))$

lemma *apply-nu-inv-mat*:

$(\text{id-blifun} - l *_R \mathcal{P}_1 d) v = \text{Bfun } (\lambda i. ((\text{nu-inv-mat } d) * v (\text{vec-lambda } v)) \$ i)$

<proof>

lemma *bounded-linear-vec-lambda*: *bounded-linear* $(\lambda x. \text{vec-lambda } (x :: 's \Rightarrow_b \text{real}))$

<proof>

lemma *bounded-linear-vec-lambda-blinfun*:
 fixes $f :: ('s \Rightarrow_b \text{real}) \Rightarrow_L ('s \Rightarrow_b \text{real})$
 shows *bounded-linear* ($\lambda v. \text{vec-lambda } (\text{apply-bfun } (\text{blinfun-apply } f$
 ($\text{bfun.Bfun } ((\$) v))))$)
 <proof>

lemma *invertible-nu-inv-max*: *invertible* (*nu-inv-mat* d)
 <proof>

end

definition *least-arg-max* $f P = (\text{LEAST } x. \text{is-arg-max } f P x)$

locale *MDP-ord* = *MDP-finite-type* $A K r l$
 for A **and**
 $K :: 's :: \{\text{finite}, \text{wellorder}\} \times 'a :: \{\text{finite}, \text{wellorder}\} \Rightarrow 's \text{ pmf}$
 and $r l$
begin

lemma *L-fin-eq-det*: $\mathcal{L} v s = (\bigsqcup a \in A s. L_a a v s)$
 <proof>

lemma *L_b-fin-eq-det*: $\mathcal{L}_b v s = (\bigsqcup a \in A s. L_a a v s)$
 <proof>

sublocale *MDP-PI-finite* $A K r l \lambda X. \text{Least } (\lambda x. x \in X)$
 <proof>

end
end

theory *Modified-Policy-Iteration*
 imports
 Policy-Iteration
 Value-Iteration
begin

3 Modified Policy Iteration

locale *MDP-MPI* = *MDP-finite-type* $A K r l + \text{MDP-act } A K r l$
 arb-act
 for A **and** $K :: 's :: \text{finite} \times 'a :: \text{finite} \Rightarrow 's \text{ pmf}$ **and** $r l \text{ arb-act}$
begin

3.1 The Advantage Function B

definition $B v s = (\bigsqcup d \in D_R. (r\text{-dec } d s + (l *_R \mathcal{P}_1 d - id\text{-blinfun } v s)))$

The function B denotes the advantage of choosing the optimal action vs. the current value estimate

lemma $B\text{-eq-}\mathcal{L}$: $B v s = \mathcal{L} v s - v s$
 $\langle proof \rangle$

B is a bounded function.

lift-definition $B_b :: ('s \Rightarrow_b real) \Rightarrow 's \Rightarrow_b real$ is B
 $\langle proof \rangle$

lemma $B_b\text{-eq-}\mathcal{L}_b$: $B_b v = \mathcal{L}_b v - v$
 $\langle proof \rangle$

lemma $\mathcal{L}_b\text{-eq-SUP-}\mathcal{L}_a$: $\mathcal{L}_b v s = (\bigsqcup a \in A s. \mathcal{L}_a a v s)$
 $\langle proof \rangle$

3.2 Optimization of the Value Function over Multiple Steps

definition $U m v s = (\bigsqcup d \in D_R. (\nu_b\text{-fin } (mk\text{-stationary } d) m + ((l *_R \mathcal{P}_1 d) \widehat{\sim} m) v s))$

U expresses the value estimate obtained by optimizing the first m steps and afterwards using the current estimate.

lemma $U\text{-zero [simp]}$: $U 0 v = v$
 $\langle proof \rangle$

lemma $U\text{-one-eq-}\mathcal{L}$: $U 1 v s = \mathcal{L} v s$
 $\langle proof \rangle$

lift-definition $U_b :: nat \Rightarrow ('s \Rightarrow_b real) \Rightarrow ('s \Rightarrow_b real)$ is U
 $\langle proof \rangle$

lemma $U_b\text{-contraction}$: $dist (U_b m v) (U_b m u) \leq l \wedge m * dist v u$
 $\langle proof \rangle$

lemma $U_b\text{-conv}$:
 $\exists! v. U_b (Suc m) v = v$
 $(\lambda n. (U_b (Suc m) \widehat{\sim} n) v) \longrightarrow (THE v. U_b (Suc m) v = v)$
 $\langle proof \rangle$

lemma $U_b\text{-convergent}$: $convergent (\lambda n. (U_b (Suc m) \widehat{\sim} n) v)$
 $\langle proof \rangle$

lemma $U_b\text{-mono}$:

assumes $v \leq u$
shows $U_b m v \leq U_b m u$
 $\langle proof \rangle$

lemma $U_b\text{-le-}\mathcal{L}_b$: $U_b m v \leq (\mathcal{L}_b \hat{\sim} m) v$
 $\langle proof \rangle$

lemma $L\text{-iter-le-}U_b$:
assumes $d \in D_R$
shows $(L d \hat{\sim} m) v \leq U_b m v$
 $\langle proof \rangle$

lemma $lim\text{-}U_b$: $lim (\lambda n. (U_b (Suc m) \hat{\sim} n) v) = \nu_b\text{-opt}$
 $\langle proof \rangle$

lemma $U_b\text{-tendsto}$: $(\lambda n. (U_b (Suc m) \hat{\sim} n) v) \longrightarrow \nu_b\text{-opt}$
 $\langle proof \rangle$

lemma $U_b\text{-fix-unique}$: $U_b (Suc m) v = v \longleftrightarrow v = \nu_b\text{-opt}$
 $\langle proof \rangle$

lemma $dist\text{-}U_b\text{-opt}$: $dist (U_b m v) \nu_b\text{-opt} \leq l \hat{\sim} m * dist v \nu_b\text{-opt}$
 $\langle proof \rangle$

3.3 Expressing a Single Step of Modified Policy Iteration

The function W equals the value computed by the Modified Policy Iteration Algorithm in a single iteration. The right hand addend in the definition describes the advantage of using the optimal action for the first m steps.

definition $W d m v = v + (\sum i < m. (l *_R \mathcal{P}_1 d) \hat{\sim} i) (B_b v)$

lemma $W\text{-eq-}L\text{-iter}$:
assumes $\nu\text{-improving } v d$
shows $W d m v = (L d \hat{\sim} m) v$
 $\langle proof \rangle$

lemma $W\text{-le-}U_b$:
assumes $v \leq u$ $\nu\text{-improving } v d$
shows $W d m v \leq U_b m u$
 $\langle proof \rangle$

lemma $W\text{-ge-}\mathcal{L}_b$:
assumes $v \leq u$ $0 \leq B_b u$ $\nu\text{-improving } u d'$

shows $\mathcal{L}_b v \leq W d' (Suc m) u$
 $\langle proof \rangle$

lemma B_b -le:

assumes ν -improving $v d$
shows $B_b v + (l *_R \mathcal{P}_1 d - id\text{-blinfun}) (u - v) \leq B_b u$
 $\langle proof \rangle$

lemma \mathcal{L}_b -W-ge:

assumes $u \leq \mathcal{L}_b u$ ν -improving $u d$
shows $W d m u \leq \mathcal{L}_b (W d m u)$
 $\langle proof \rangle$

3.4 Computing the Bellman Operator over Multiple Steps

definition $L\text{-pow} :: ('s \Rightarrow_b \text{real}) \Rightarrow ('s \Rightarrow 'a) \Rightarrow \text{nat} \Rightarrow ('s \Rightarrow_b \text{real})$
where

$L\text{-pow } v d m = (L (mk\text{-dec-det } d) \overset{\sim}{\sim} Suc m) v$

lemma $sum\text{-telescope}'$: $(\sum_{i \leq k}. f (Suc i) - f i) = f (Suc k) - (f 0)$
 $:: 'c :: ab\text{-group-add}$
 $\langle proof \rangle$

lemma $L\text{-pow-eq}$:

assumes ν -improving $v (mk\text{-dec-det } d)$
shows $L\text{-pow } v d m = v + (\sum_{i \leq m}. ((l *_R \mathcal{P}_1 (mk\text{-dec-det } d)) \overset{\sim}{\sim} i))$
 $(B_b v)$
 $\langle proof \rangle$

lemma $L\text{-pow-eq-W}$:

assumes $d \in D_D$
shows $L\text{-pow } v (policy\text{-improvement } d v) m = W (mk\text{-dec-det } (policy\text{-improvement } d v)) (Suc m) v$
 $\langle proof \rangle$

lemma $L\text{-pow-}\mathcal{L}_b\text{-mono-inv}$:

assumes $d \in D_D$ $v \leq \mathcal{L}_b v$
shows $L\text{-pow } v (policy\text{-improvement } d v) m \leq \mathcal{L}_b (L\text{-pow } v (policy\text{-improvement } d v) m)$
 $\langle proof \rangle$

3.5 The Modified Policy Iteration Algorithm

context

fixes $d0 :: 's \Rightarrow 'a$
fixes $v0 :: 's \Rightarrow_b \text{real}$
fixes $m :: \text{nat} \Rightarrow ('s \Rightarrow_b \text{real}) \Rightarrow \text{nat}$

assumes $d0$: $d0 \in D_D$
begin

We first define a function that executes the algorithm for n steps.

fun $mpi :: nat \Rightarrow (('s \Rightarrow 'a) \times ('s \Rightarrow_b real))$ **where**
 $mpi\ 0 = (policy-improvement\ d0\ v0, v0)$ |
 $mpi\ (Suc\ n) =$
 $(let\ (d, v) = mpi\ n; v' = L-pow\ v\ d\ (m\ n\ v)\ in$
 $(policy-improvement\ d\ v', v'))$

definition $mpi-val\ n = snd\ (mpi\ n)$

definition $mpi-pol\ n = fst\ (mpi\ n)$

lemma $mpi-pol-zero[simp]$: $mpi-pol\ 0 = policy-improvement\ d0\ v0$
 $\langle proof \rangle$

lemma $mpi-pol-Suc$: $mpi-pol\ (Suc\ n) = policy-improvement\ (mpi-pol\ n)\ (mpi-val\ (Suc\ n))$
 $\langle proof \rangle$

lemma $mpi-pol-is-dec-det$: $mpi-pol\ n \in D_D$
 $\langle proof \rangle$

lemma ν -*improving-mpi-pol*: ν -*improving* $(mpi-val\ n)\ (mk-dec-det\ (mpi-pol\ n))$
 $\langle proof \rangle$

lemma $mpi-val-zero[simp]$: $mpi-val\ 0 = v0$
 $\langle proof \rangle$

lemma $mpi-val-Suc$: $mpi-val\ (Suc\ n) = L-pow\ (mpi-val\ n)\ (mpi-pol\ n)\ (m\ n\ (mpi-val\ n))$
 $\langle proof \rangle$

lemma $mpi-val-eq$: $mpi-val\ (Suc\ n) =$
 $mpi-val\ n + (\sum\ i \leq m\ n\ (mpi-val\ n). (l *_R\ \mathcal{P}_1\ (mk-dec-det\ (mpi-pol\ n))) \wedge i)\ (B_b\ (mpi-val\ n))$
 $\langle proof \rangle$

Value Iteration is a special case of MPI where $\forall n\ v. m\ n\ v = 0$.

lemma mpi -*includes-value-it*:

assumes $\forall n\ v. m\ n\ v = 0$

shows $mpi-val\ (Suc\ n) = \mathcal{L}_b\ (mpi-val\ n)$

$\langle proof \rangle$

3.6 Convergence Proof

We define the sequence w as an upper bound for the values of MPI.

fun w **where**

$w\ 0 = v0 \mid$

$w\ (Suc\ n) = U_b\ (Suc\ (m\ n\ (mpi-val\ n)))\ (w\ n)$

lemma $dist-\nu_b-opt$: $dist\ (w\ (Suc\ n))\ \nu_b-opt \leq l * dist\ (w\ n)\ \nu_b-opt$
 $\langle proof \rangle$

lemma $dist-\nu_b-opt-n$: $dist\ (w\ n)\ \nu_b-opt \leq l^n * dist\ v0\ \nu_b-opt$
 $\langle proof \rangle$

lemma $w-conv$: $w \longrightarrow \nu_b-opt$
 $\langle proof \rangle$

MPI converges monotonically to the optimal value from below. The iterates are sandwiched between \mathcal{L}_b from below and U_b from above.

theorem $mpi-conv$:

assumes $v0 \leq \mathcal{L}_b\ v0$

shows $mpi-val \longrightarrow \nu_b-opt$ **and** $\bigwedge n. mpi-val\ n \leq mpi-val\ (Suc\ n)$

$\langle proof \rangle$

3.7 ϵ -Optimality

This gives an upper bound on the error of MPI.

lemma $mpi-pol-eps-opt$:

assumes $2 * l * dist\ (mpi-val\ n)\ (\mathcal{L}_b\ (mpi-val\ n)) < eps * (1 - l)$

$eps > 0$

shows $dist\ (\nu_b\ (mk-stationary-det\ (mpi-pol\ n)))\ (\mathcal{L}_b\ (mpi-val\ n)) \leq$

$eps / 2$

$\langle proof \rangle$

lemma $mpi-pol-opt$:

assumes $2 * l * dist\ (mpi-val\ n)\ (\mathcal{L}_b\ (mpi-val\ n)) < eps * (1 - l)$

$eps > 0$

shows $dist\ (\nu_b\ (mk-stationary-det\ (mpi-pol\ n)))\ (\nu_b-opt) < eps$

$\langle proof \rangle$

lemma $mpi-val-term-ex$:

assumes $v0 \leq \mathcal{L}_b\ v0$ $eps > 0$

shows $\exists n. 2 * l * dist\ (mpi-val\ n)\ (\mathcal{L}_b\ (mpi-val\ n)) < eps * (1 - l)$

$\langle proof \rangle$

end

3.8 Unbounded MPI

context

fixes $eps\ \delta :: real$ **and** $M :: nat$

begin

function (*domintros*) *mpi-algo* **where** *mpi-algo* $d\ v\ m =$ (
 if $2 * l * \text{dist } v (\mathcal{L}_b\ v) < \text{eps} * (1 - l)$
 then (*policy-improvement* $d\ v, v$)
 else *mpi-algo* (*policy-improvement* $d\ v$) (*L-pow* v (*policy-improvement*
 $d\ v$) ($m\ 0\ v$)) ($\lambda n. m\ (\text{Suc } n)$))
 <proof>)

We define a tailrecursive version of *mpi* which more closely resembles *mpi-algo*.

fun *mpi'* **where**
 mpi' $d\ v\ 0\ m =$ (*policy-improvement* $d\ v, v$) |
 mpi' $d\ v\ (\text{Suc } n)\ m =$ (
 let $d' = \text{policy-improvement } d\ v; v' = \text{L-pow } v\ d' (m\ 0\ v)$ *in* *mpi'* $d'\ v'\ n$
 ($\lambda n. m\ (\text{Suc } n)$))

lemma *mpi-Suc'*:
 assumes $d \in D_D$
 shows $\text{mpi } d\ v\ m\ (\text{Suc } n) = \text{mpi } (\text{policy-improvement } d\ v) (\text{L-pow } v$
 (*policy-improvement* $d\ v$) ($m\ 0\ v$)) ($\lambda a. m\ (\text{Suc } a)$) n
 <proof>

lemma
 assumes $d \in D_D$
 shows $\text{mpi } d\ v\ m\ n = \text{mpi}'\ d\ v\ n\ m$
 <proof>

lemma *termination-mpi-algo*:
 assumes $\text{eps} > 0\ d \in D_D\ v \leq \mathcal{L}_b\ v$
 shows *mpi-algo-dom* (d, v, m)
 <proof>

abbreviation *mpi-alg-rec* $d\ v\ m \equiv$
 (*if* $2 * l * \text{dist } v (\mathcal{L}_b\ v) < \text{eps} * (1 - l)$ *then* (*policy-improvement*
 $d\ v, v$)
 else *mpi-algo* (*policy-improvement* $d\ v$) (*L-pow* v (*policy-improvement*
 $d\ v$) ($m\ 0\ v$))
 ($\lambda n. m\ (\text{Suc } n)$))

lemma *mpi-algo-def'*:
 assumes $d \in D_D\ v \leq \mathcal{L}_b\ v\ \text{eps} > 0$
 shows $\text{mpi-algo } d\ v\ m = \text{mpi-alg-rec } d\ v\ m$
 <proof>

lemma *mpi-algo-eq-mpi*:
 assumes $d \in D_D\ v \leq \mathcal{L}_b\ v\ \text{eps} > 0$
 shows $\text{mpi-algo } d\ v\ m = \text{mpi } d\ v\ m$ (*LEAST* $n. 2 * l * \text{dist } (\text{mpi-val}$
 $d\ v\ m\ n) (\mathcal{L}_b\ (\text{mpi-val } d\ v\ m\ n)) < \text{eps} * (1 - l)$)
 <proof>

lemma *mpi-algo-opt*:
assumes $v0 \leq \mathcal{L}_b$ $v0$ $eps > 0$ $d \in D_D$
shows $dist (\nu_b (mk-stationary-det (fst (mpi-algo d v0 m)))) \nu_b-opt < eps$
 $\langle proof \rangle$

end

3.9 Initial Value Estimate *v0-mpi*

We define an initial estimate of the value function for which Modified Policy Iteration always terminates.

abbreviation $r-min \equiv (\prod s' a. r (s', a))$

definition $v0-mpi s = r-min / (1 - l)$

lift-definition $v0-mpi_b :: 's \Rightarrow_b real$ **is** $v0-mpi$
 $\langle proof \rangle$

lemma $v0-mpi_b-le-\mathcal{L}_b$: $v0-mpi_b \leq \mathcal{L}_b$ $v0-mpi_b$
 $\langle proof \rangle$

3.10 An Instance of Modified Policy Iteration with a Valid Conservative Initial Value Estimate

definition $mpi-user eps m = ($
 $if eps \leq 0$ *then undefined else* $mpi-alg-rec eps (\lambda x. arb-act (A x))$
 $v0-mpi_b m)$

lemma *mpi-user-eq*:
assumes $eps > 0$
shows $mpi-user eps = mpi-alg-rec eps (\lambda x. arb-act (A x)) v0-mpi_b$
 $\langle proof \rangle$

lemma *mpi-user-opt*:
assumes $eps > 0$
shows $dist (\nu_b (mk-stationary-det (fst (mpi-user eps n)))) \nu_b-opt < eps$
 $\langle proof \rangle$

end

end
theory *Matrix-Util*
imports *HOL-Analysis.Analysis*
begin

4 Matrices

proposition *scalar-matrix-assoc'*:
fixes $C :: ('b::\text{real-algebra-1})^m{}^n$
shows $k *_R (C ** D) = C ** (k *_R D)$
 $\langle\text{proof}\rangle$

4.1 Nonnegative Matrices

lemma *nonneg-matrix-nonneg* [*dest*]: $0 \leq m \implies 0 \leq m \$ i \$ j$
 $\langle\text{proof}\rangle$

lemma *matrix-mult-mono*:
assumes $0 \leq E \ 0 \leq C \ (E :: \text{real}^c{}^c) \leq B \ C \leq D$
shows $E ** C \leq B ** D$
 $\langle\text{proof}\rangle$

lemma *nonneg-matrix-mult*: $0 \leq (C :: ('b::\{\text{field}, \text{ordered-ring}\})^{\wedge}{}^{\wedge})$
 $\implies 0 \leq D \implies 0 \leq C ** D$
 $\langle\text{proof}\rangle$

lemma *zero-le-mat-iff* [*simp*]: $0 \leq \text{mat } (x :: 'c :: \{\text{zero}, \text{order}\}) \longleftrightarrow$
 $0 \leq x$
 $\langle\text{proof}\rangle$

lemma *nonneg-mat-ge-zero*: $0 \leq Q \implies 0 \leq v \implies 0 \leq Q *v \ (v ::$
 $\text{real}^c)$
 $\langle\text{proof}\rangle$

lemma *nonneg-mat-mono*: $0 \leq Q \implies u \leq v \implies Q *v \ u \leq Q *v \ (v$
 $:: \text{real}^c)$
 $\langle\text{proof}\rangle$

lemma *nonneg-mult-imp-nonneg-mat*:
assumes $\bigwedge v. v \geq 0 \implies X *v \ v \geq 0$
shows $X \geq (0 :: \text{real}^{\wedge}{}^{\wedge})$
 $\langle\text{proof}\rangle$

lemma *nonneg-mat-iff*:
 $(X \geq (0 :: \text{real}^{\wedge}{}^{\wedge})) \longleftrightarrow (\forall v. v \geq 0 \implies X *v \ v \geq 0)$
 $\langle\text{proof}\rangle$

lemma *mat-le-iff*: $(X \leq Y) \longleftrightarrow (\forall x \geq 0. (X :: \text{real}^{\wedge}{}^{\wedge}) *v \ x \leq Y *v$
 $x)$
 $\langle\text{proof}\rangle$

4.2 Matrix Powers

primrec *matpow* :: $'a::\text{semiring-1}^n{}^n \Rightarrow \text{nat} \Rightarrow 'a^{\wedge}{}^n$ **where**
matpow-0: $\text{matpow } A \ 0 = \text{mat } 1 \ |$

matpow-Suc: $\text{matpow } A \text{ (Suc } n) = (\text{matpow } A \text{ } n) ** A$

lemma *nonneg-matpow*: $0 \leq X \implies 0 \leq \text{matpow } (X :: \text{real } ^\wedge - \wedge -) \ i$
 $\langle \text{proof} \rangle$

lemma *matpow-mono*: $0 \leq C \implies C \leq D \implies \text{matpow } (C :: \text{real } ^\wedge - \wedge -)$
 $n \leq \text{matpow } D \ n$
 $\langle \text{proof} \rangle$

lemma *matpow-scaleR*: $\text{matpow } (c *_{\mathbb{R}} (X :: 'b :: \text{real-algebra-1 } ^\wedge - \wedge -))$
 $n = (c \wedge n) *_{\mathbb{R}} (\text{matpow } X) \ n$
 $\langle \text{proof} \rangle$

lemma *matrix-vector-mult-code'*: $(X *_{\mathbb{V}} x) \$ i = (\sum_{j \in \text{UNIV}} X \$ i \ j$
 $\$ j * x \$ j)$
 $\langle \text{proof} \rangle$

lemma *matrix-vector-mult-mono*: $(0 :: \text{real } ^\wedge - \wedge -) \leq X \implies 0 \leq v \implies X$
 $\leq Y \implies X *_{\mathbb{V}} v \leq Y *_{\mathbb{V}} v$
 $\langle \text{proof} \rangle$

4.3 Triangular Matrices

definition *lower-triangular-mat* $X \longleftrightarrow (\forall i \ j. (i :: 'b :: \{\text{finite}, \text{linorder}\})$
 $< j \longrightarrow X \$ i \$ j = 0)$

definition *strict-lower-triangular-mat* $X \longleftrightarrow (\forall i \ j. (i :: 'b :: \{\text{finite},$
 $\text{linorder}\}) \leq j \longrightarrow X \$ i \$ j = 0)$

definition *upper-triangular-mat* $X \longleftrightarrow (\forall i \ j. j < i \longrightarrow X \$ i \$ j =$
 $0)$

lemma *stlI*: *strict-lower-triangular-mat* $X \implies$ *lower-triangular-mat* X
 $\langle \text{proof} \rangle$

lemma *lower-triangular-mat-mat*: *lower-triangular-mat* $(\text{mat } x)$
 $\langle \text{proof} \rangle$

lemma *lower-triangular-mult*:
assumes *lower-triangular-mat* X *lower-triangular-mat* Y
shows *lower-triangular-mat* $(X ** Y)$
 $\langle \text{proof} \rangle$

lemma *lower-triangular-pow*:
assumes *lower-triangular-mat* X
shows *lower-triangular-mat* $(\text{matpow } X \ i)$
 $\langle \text{proof} \rangle$

lemma *lower-triangular-suminf*:
assumes $\bigwedge i. \text{lower-triangular-mat } (f \ i) \text{ summable } (f \ :: \ \text{nat} \Rightarrow \text{'b::real-normed-vector } \hat{\ } \hat{\ })$
shows $\text{lower-triangular-mat } (\sum i. f \ i)$
 $\langle \text{proof} \rangle$

lemma *lower-triangular-pow-eq*:
assumes $\text{lower-triangular-mat } X \ \text{lower-triangular-mat } Y \ \bigwedge s'. s' \leq s \Rightarrow \text{row } s' \ X = \text{row } s' \ Y \ s' \leq s$
shows $\text{row } s' \ (\text{matpow } X \ i) = \text{row } s' \ (\text{matpow } Y \ i)$
 $\langle \text{proof} \rangle$

lemma *lower-triangular-mat-mult*:
assumes $\text{lower-triangular-mat } M \ \bigwedge i. i \leq j \Rightarrow v \ \$ \ i = v' \ \$ \ i$
shows $(M \ *v \ v) \ \$ \ j = (M \ *v \ v') \ \$ \ j$
 $\langle \text{proof} \rangle$

4.4 Inverses

lemma *matrix-inv*:
assumes *invertible* M
shows *matrix-inv-left*: $\text{matrix-inv } M \ ** \ M = \text{mat } 1$
and *matrix-inv-right*: $M \ ** \ \text{matrix-inv } M = \text{mat } 1$
 $\langle \text{proof} \rangle$

lemma *matrix-inv-unique*:
fixes $A::\text{'a::}\{\text{semiring-1}\}^{\sim n} \sim n$
assumes *AB*: $A \ ** \ B = \text{mat } 1$ **and** *BA*: $B \ ** \ A = \text{mat } 1$
shows $\text{matrix-inv } A = B$
 $\langle \text{proof} \rangle$

end

theory *Blinfun-Matrix*

imports

MDP-Rewards.Blinfun-Util

Matrix-Util

begin

5 Bounded Linear Functions and Matrices

definition *blinfun-to-matrix* $(f \ :: \ (\text{'b::finite} \Rightarrow_b \ \text{real}) \Rightarrow_L \ (\text{'c::finite} \Rightarrow_b \ -)) =$
 $\text{matrix } (\lambda v. (\chi \ j. f \ (\text{Bfun } ((\$) \ v)) \ j))$

definition *matrix-to-blinfun* $X = \text{Blinfun } (\lambda v. \text{Bfun } (\lambda i. (X \ *v \ (\chi \ i. (\text{apply-bfun } v \ i))) \ \$ \ i))$

lemma *plus-vec-eq*: $(\chi \ i. f \ i + g \ i) = (\chi \ i. f \ i) + (\chi \ i. g \ i)$

$\langle \text{proof} \rangle$

lemma *matrix-to-blinfun-mult*: $\text{matrix-to-blinfun } m \ (v :: 'c::\text{finite} \Rightarrow_b \text{real}) \ i = (m *v (\chi \ i. \ v \ i)) \ \$ \ i$
 $\langle \text{proof} \rangle$

lemma *blinfun-to-matrix-mult*: $(\text{blinfun-to-matrix } f *v (\chi \ i. \ \text{apply-bfun } v \ i)) \ \$ \ i = f \ v \ i$
 $\langle \text{proof} \rangle$

lemma *blinfun-to-matrix-mult'*: $(\text{blinfun-to-matrix } f *v \ v) \ \$ \ i = f \ (\text{Bfun } (\lambda i. \ v \ \$ \ i)) \ i$
 $\langle \text{proof} \rangle$

lemma *blinfun-to-matrix-mult''*: $(\text{blinfun-to-matrix } f *v \ v) = (\chi \ i. \ f \ (\text{Bfun } (\lambda i. \ v \ \$ \ i)) \ i)$
 $\langle \text{proof} \rangle$

lemma *matrix-to-blinfun-inv*: $\text{matrix-to-blinfun } (\text{blinfun-to-matrix } f) = f$
 $\langle \text{proof} \rangle$

lemma *blinfun-to-matrix-add*: $\text{blinfun-to-matrix } (f + g) = \text{blinfun-to-matrix } f + \text{blinfun-to-matrix } g$
 $\langle \text{proof} \rangle$

lemma *blinfun-to-matrix-diff*: $\text{blinfun-to-matrix } (f - g) = \text{blinfun-to-matrix } f - \text{blinfun-to-matrix } g$
 $\langle \text{proof} \rangle$

lemma *blinfun-to-matrix-scaleR*: $\text{blinfun-to-matrix } (c *R \ f) = c *R \ \text{blinfun-to-matrix } f$
 $\langle \text{proof} \rangle$

lemma *matrix-to-blinfun-add*:
 $\text{matrix-to-blinfun } ((f :: \text{real}^{\wedge} \text{-} \wedge) + g) = \text{matrix-to-blinfun } f + \text{matrix-to-blinfun } g$
 $\langle \text{proof} \rangle$

lemma *matrix-to-blinfun-diff*:
 $\text{matrix-to-blinfun } ((f :: \text{real}^{\wedge} \text{-} \wedge) - g) = \text{matrix-to-blinfun } f - \text{matrix-to-blinfun } g$
 $\langle \text{proof} \rangle$

lemma *matrix-to-blinfun-scaleR*:
 $\text{matrix-to-blinfun } (c *R \ (f :: \text{real}^{\wedge} \text{-} \wedge)) = c *R \ \text{matrix-to-blinfun } f$
 $\langle \text{proof} \rangle$

lemma *matrix-to-blinfun-comp*: $\text{matrix-to-blinfun } ((m :: \text{real}^{\wedge} \text{-} \wedge) **$

$n) = (\text{matrix-to-blinfun } m) \circ_L (\text{matrix-to-blinfun } n)$
 ⟨proof⟩

lemma *blinfun-to-matrix-comp*: $\text{blinfun-to-matrix } (f \circ_L g) = (\text{blinfun-to-matrix } f) ** (\text{blinfun-to-matrix } g)$
 ⟨proof⟩

lemma *blinfun-to-matrix-id*: $\text{blinfun-to-matrix } \text{id-blinfun} = \text{mat } 1$
 ⟨proof⟩

lemma *matrix-to-blinfun-id*: $\text{matrix-to-blinfun } (\text{mat } 1 :: (\text{real } \hat{-} \hat{-})) = \text{id-blinfun}$
 ⟨proof⟩

lemma *matrix-to-blinfun-inv_L*:
assumes *invertible* m
shows $\text{matrix-to-blinfun } (\text{matrix-inv } (m :: \text{real } \hat{-} \hat{-})) = \text{inv}_L (\text{matrix-to-blinfun } m)$
 $\text{invertible}_L (\text{matrix-to-blinfun } m)$
 ⟨proof⟩

lemma *blinfun-to-matrix-inverse*:
assumes *invertible_L* X
shows *invertible* $(\text{blinfun-to-matrix } (X :: ('b :: \text{finite} \Rightarrow_b \text{real}) \Rightarrow_L 'c :: \text{finite} \Rightarrow_b \text{real}))$
 $\text{blinfun-to-matrix } (\text{inv}_L X) = \text{matrix-inv } (\text{blinfun-to-matrix } X)$
 ⟨proof⟩

lemma *blinfun-to-matrix-inv[simp]*: $\text{blinfun-to-matrix } (\text{matrix-to-blinfun } f) = f$
 ⟨proof⟩

lemma *invertible-invertible_L-I*: $\text{invertible } (\text{blinfun-to-matrix } f) \Longrightarrow \text{invertible}_L f$
 $\text{invertible}_L (\text{matrix-to-blinfun } X) \Longrightarrow \text{invertible } (X :: \text{real } \hat{-} \hat{-})$
 ⟨proof⟩

lemma *bounded-linear-blinfun-to-matrix*: $\text{bounded-linear } (\text{blinfun-to-matrix } :: ('a \Rightarrow_b \text{real}) \Rightarrow_L ('b \Rightarrow_b \text{real}) \Rightarrow \text{real } \hat{-} 'a \hat{-} 'b)$
 ⟨proof⟩

lemma *summable-blinfun-to-matrix*:
assumes *summable* $(f :: \text{nat} \Rightarrow ('c :: \text{finite} \Rightarrow_b -) \Rightarrow_L ('c \Rightarrow_b -))$
shows *summable* $(\lambda i. \text{blinfun-to-matrix } (f i))$
 ⟨proof⟩

abbreviation *nonneg-blinfun* $Q \equiv 0 \leq (\text{blinfun-to-matrix } Q)$

lemma *nonneg-blinfun-mono*: $\text{nonneg-blinfun } Q \implies u \leq v \implies Q u \leq Q v$
 ⟨proof⟩

lemma *nonneg-blinfun-nonneg*: $\text{nonneg-blinfun } Q \implies 0 \leq v \implies 0 \leq Q v$
 ⟨proof⟩

lemma *nonneg-id-blinfun*: $\text{nonneg-blinfun id-blinfun}$
 ⟨proof⟩

lemma *norm-nonneg-blinfun-one*:
assumes $0 \leq \text{blinfun-to-matrix } X$
shows $\text{norm } X = \text{norm } (\text{blinfun-apply } X 1)$
 ⟨proof⟩

lemma *matrix-le-norm-mono*:
assumes $0 \leq (\text{blinfun-to-matrix } C)$
and $(\text{blinfun-to-matrix } C) \leq (\text{blinfun-to-matrix } D)$
shows $\text{norm } C \leq \text{norm } D$
 ⟨proof⟩

lemma *blinfun-to-matrix-matpow*: $\text{blinfun-to-matrix } (X \hat{\sim} i) = \text{matpow } (\text{blinfun-to-matrix } X) i$
 ⟨proof⟩

lemma *nonneg-blinfun-iff*: $\text{nonneg-blinfun } X \iff (\forall v \geq 0. X v \geq 0)$
 ⟨proof⟩

lemma *blinfun-apply-mono*: $(0::\text{real}^{\hat{\sim}}\hat{\sim}) \leq \text{blinfun-to-matrix } X \implies 0 \leq v \implies \text{blinfun-to-matrix } X \leq \text{blinfun-to-matrix } Y \implies X v \leq Y v$
 ⟨proof⟩

end

theory *Splitting-Methods*

imports

Blinfun-Matrix

Value-Iteration

Policy-Iteration

begin

6 Value Iteration using Splitting Methods

6.1 Regular Splittings for Matrices and Bounded Linear Functions

definition *is-splitting-mat* $X Q R \iff$

$$X = Q - R \wedge \text{invertible } Q \wedge 0 \leq \text{matrix-inv } Q \wedge 0 \leq R$$

definition *is-splitting-blin* $X Q R \longleftrightarrow \text{is-splitting-mat } (\text{blinfun-to-matrix } X) (\text{blinfun-to-matrix } Q) (\text{blinfun-to-matrix } R)$

lemma *is-splitting-blin-def'*: *is-splitting-blin* $X Q R \longleftrightarrow$

$X = Q - R \wedge \text{invertible}_L Q \wedge \text{nonneg-blinfun } (\text{inv}_L Q) \wedge \text{nonneg-blinfun } R$
 ⟨proof⟩

lemma *is-splitting-blinD[dest]*:

assumes *is-splitting-blin* $X Q R$
shows $X = Q - R \wedge \text{invertible}_L Q \wedge \text{nonneg-blinfun } (\text{inv}_L Q) \wedge \text{nonneg-blinfun } R$
 ⟨proof⟩

6.2 Splitting Methods for MDPs

locale *MDP-QR* = *MDP-finite-type* $A K r l$

for $A :: 's :: \text{finite} \Rightarrow ('a :: \text{finite}) \text{ set}$
and $K :: ('s \times 'a) \Rightarrow 's \text{ pmf}$
and $r l +$
fixes $Q :: ('s \Rightarrow 'a) \Rightarrow ('s \Rightarrow_b \text{real}) \Rightarrow_L ('s \Rightarrow_b \text{real})$
fixes $R :: ('s \Rightarrow 'a) \Rightarrow ('s \Rightarrow_b \text{real}) \Rightarrow_L ('s \Rightarrow_b \text{real})$
assumes *is-splitting*: $\bigwedge d. d \in D_D \Longrightarrow \text{is-splitting-blin } (\text{id-blinfun } - l *_R \mathcal{P}_1 (\text{mk-dec-det } d)) (Q d) (R d)$
assumes *QR-contraction*: $(\bigsqcup_{d \in D_D} \text{norm } (\text{inv}_L (Q d) o_L R d)) < 1$
assumes *arg-max-ex-split*: $\exists d. \forall s. \text{is-arg-max } (\lambda d. \text{inv}_L (Q d) (r\text{-dec}_b (\text{mk-dec-det } d) + R d v) s) (\lambda d. d \in D_D) d$
begin

lemma *inv-Q-mono*: $d \in D_D \Longrightarrow u \leq v \Longrightarrow (\text{inv}_L (Q d)) u \leq (\text{inv}_L (Q d)) v$
 ⟨proof⟩

lemma *splitting-eq*: $d \in D_D \Longrightarrow Q d - R d = (\text{id-blinfun } - l *_R \mathcal{P}_1 (\text{mk-dec-det } d))$
 ⟨proof⟩

lemma *Q-nonneg*: $d \in D_D \Longrightarrow 0 \leq v \Longrightarrow 0 \leq \text{inv}_L (Q d) v$
 ⟨proof⟩

lemma *Q-invertible*: $d \in D_D \Longrightarrow \text{invertible}_L (Q d)$
 ⟨proof⟩

lemma *R-nonneg*: $d \in D_D \Longrightarrow 0 \leq v \Longrightarrow 0 \leq R d v$
 ⟨proof⟩

lemma *R-mono*: $d \in D_D \implies u \leq v \implies (R d) u \leq (R d) v$
 ⟨proof⟩

lemma *QR-nonneg*: $d \in D_D \implies 0 \leq v \implies 0 \leq (inv_L (Q d) o_L R d) v$
 ⟨proof⟩

lemma *QR-mono*: $d \in D_D \implies u \leq v \implies (inv_L (Q d) o_L R d) u \leq (inv_L (Q d) o_L R d) v$
 ⟨proof⟩

lemma *norm-QR-less-one*: $d \in D_D \implies norm (inv_L (Q d) o_L R d) < 1$
 ⟨proof⟩

lemma *splitting*: $d \in D_D \implies id\text{-blinfun} - l *_R \mathcal{P}_1 (mk\text{-dec-det } d) = Q d - R d$
 ⟨proof⟩

6.3 Discount Factor *QR-disc*

abbreviation *QR-disc* $\equiv (\bigsqcup d \in D_D. norm (inv_L (Q d) o_L R d))$

lemma *QR-le-QR-disc*: $d \in D_D \implies norm (inv_L (Q d) o_L (R d)) \leq QR\text{-disc}$
 ⟨proof⟩

lemma *a-nonneg*: $0 \leq QR\text{-disc}$
 ⟨proof⟩

6.4 Bellman-Operator

abbreviation *L-split* $d v \equiv inv_L (Q d) (r\text{-dec}_b (mk\text{-dec-det } d) + R d v)$

definition *L-split* $v s = (\bigsqcup d \in D_D. L\text{-split } d v s)$

lemma *L-split-bfun-aux*:

assumes $d \in D_D$

shows $norm (L\text{-split } d v) \leq (\bigsqcup d \in D_D. norm (inv_L (Q d))) * r_M + norm v$

⟨proof⟩

lift-definition *L_b-split* :: $('s \Rightarrow_b real) \Rightarrow ('s \Rightarrow_b real)$ **is** *L-split*
 ⟨proof⟩

lemma *L_b-split-def'*: $L_b\text{-split } v s = (\bigsqcup d \in D_D. L\text{-split } d v s)$
 ⟨proof⟩

lemma \mathcal{L}_b -split-contraction: $\text{dist } (\mathcal{L}_b\text{-split } v) (\mathcal{L}_b\text{-split } u) \leq \text{QR-disc} * \text{dist } v u$

$\langle \text{proof} \rangle$

lemma \mathcal{L}_b -lim:

$\exists ! v. \mathcal{L}_b\text{-split } v = v$

$(\lambda n. (\mathcal{L}_b\text{-split } \widetilde{\sim} n) v) \longrightarrow (\text{THE } v. \mathcal{L}_b\text{-split } v = v)$

$\langle \text{proof} \rangle$

lemma \mathcal{L}_b -split-tendsto-opt: $(\lambda n. (\mathcal{L}_b\text{-split } \widetilde{\sim} n) v) \longrightarrow \nu_b\text{-opt}$

$\langle \text{proof} \rangle$

lemma \mathcal{L}_b -split-fix[simp]: $\mathcal{L}_b\text{-split } \nu_b\text{-opt} = \nu_b\text{-opt}$

$\langle \text{proof} \rangle$

lemma $\text{dist-}\mathcal{L}_b\text{-split-opt-eps}$:

assumes $\text{eps} > 0 \ 2 * \text{QR-disc} * \text{dist } v (\mathcal{L}_b\text{-split } v) < \text{eps} * (1 - \text{QR-disc})$

shows $\text{dist } (\mathcal{L}_b\text{-split } v) \nu_b\text{-opt} < \text{eps} / 2$

$\langle \text{proof} \rangle$

lemma L -split-fix:

assumes $d \in D_D$

shows $L\text{-split } d (\nu_b (\text{mk-stationary-det } d)) = \nu_b (\text{mk-stationary-det } d)$

$\langle \text{proof} \rangle$

lemma L -split-contraction:

assumes $d \in D_D$

shows $\text{dist } (L\text{-split } d v) (L\text{-split } d u) \leq \text{QR-disc} * \text{dist } v u$

$\langle \text{proof} \rangle$

lemma $\text{find-policy-QR-error-bound}$:

assumes $\text{eps} > 0 \ 2 * \text{QR-disc} * \text{dist } v (\mathcal{L}_b\text{-split } v) < \text{eps} * (1 - \text{QR-disc})$

assumes $\text{am}: \bigwedge s. \text{is-arg-max } (\lambda d. L\text{-split } d (\mathcal{L}_b\text{-split } v) s) (\lambda d. d \in D_D) d$

shows $\text{dist } (\nu_b (\text{mk-stationary-det } d)) \nu_b\text{-opt} < \text{eps}$

$\langle \text{proof} \rangle$

end

context MDP-ord **begin**

lemma $\text{inv-one-sub-Q}'$:

fixes $Q :: 'c :: \text{banach} \Rightarrow_L 'c$

assumes $\text{onorm-le}: \text{norm } (\text{id-blifun} - Q) < 1$

shows $\text{inv}_L Q = (\sum i. (\text{id-blifun} - Q) \widetilde{\sim} i)$

$\langle \text{proof} \rangle$

An important theorem: allows to compare the rate of convergence

for different splittings

lemma *norm-splitting-le*:

assumes *is-splitting-blin* ($id\text{-blinfun} - l *_R \mathcal{P}_1 d$) $Q1 R1$
and *is-splitting-blin* ($id\text{-blinfun} - l *_R \mathcal{P}_1 d$) $Q2 R2$
and ($blinfun\text{-to-matrix } R2 \leq blinfun\text{-to-matrix } R1$)
and ($blinfun\text{-to-matrix } R1 \leq blinfun\text{-to-matrix } (l *_R \mathcal{P}_1 d)$)
shows $norm (inv_L Q2 o_L R2) \leq norm (inv_L Q1 o_L R1)$

<proof>

6.5 Gauss Seidel Splitting

6.5.1 Definition of Upper and Lower Triangular Matrices

definition *P-dec* $d \equiv blinfun\text{-to-matrix } (\mathcal{P}_1 (mk\text{-dec-det } d))$

definition *P-upper* $d \equiv (\chi \ i \ j. \text{ if } i \leq j \text{ then } P\text{-dec } d \ \$ \ i \ \$ \ j \text{ else } 0)$

definition *P-lower* $d \equiv (\chi \ i \ j. \text{ if } j < i \text{ then } P\text{-dec } d \ \$ \ i \ \$ \ j \text{ else } 0)$

definition $\mathcal{P}_U \ d = matrix\text{-to-blinfun } (P\text{-upper } d)$

definition $\mathcal{P}_L \ d = matrix\text{-to-blinfun } (P\text{-lower } d)$

lemma *P-dec-elem*: $P\text{-dec } d \ \$ \ i \ \$ \ j = pmf (K (i, d i)) j$

<proof>

lemma *nonneg- \mathcal{P}_U* : $nonneg\text{-blinfun } (\mathcal{P}_U \ d)$

<proof>

lemma *nonneg-P-dec*: $0 \leq P\text{-dec } d$

<proof>

lemma *nonneg-P-upper*: $0 \leq P\text{-upper } d$

<proof>

lemma *nonneg-P-lower*: $0 \leq P\text{-lower } d$

<proof>

lemma *nonneg- \mathcal{P}_L* : $nonneg\text{-blinfun } (\mathcal{P}_L \ d)$

<proof>

lemma *nonneg- \mathcal{P}_1* : $nonneg\text{-blinfun } (\mathcal{P}_1 \ d)$

<proof>

lemma *norm- \mathcal{P}_L -le*: $norm (\mathcal{P}_L \ d) \leq norm (\mathcal{P}_1 (mk\text{-dec-det } d))$

<proof>

lemma *norm- \mathcal{P}_L -le-one*: $norm (\mathcal{P}_L \ d) \leq 1$

<proof>

lemma *norm- \mathcal{P}_L -less-one*: $norm (l *_R \mathcal{P}_L \ d) < 1$

<proof>

lemma $\mathcal{P}_L\text{-le-}\mathcal{P}_1$: $0 \leq v \implies \mathcal{P}_L d v \leq \mathcal{P}_1 (mk\text{-dec-det } d) v$
<proof>

lemma $\mathcal{P}_U\text{-le-}\mathcal{P}_1$: $0 \leq v \implies \mathcal{P}_U d v \leq \mathcal{P}_1 (mk\text{-dec-det } d) v$
<proof>

lemma *row-P-upper-indep*: $d s = d' s \implies \text{row } s (P\text{-upper } d) = \text{row } s (P\text{-upper } d')$
<proof>

lemma *row-P-lower-indep*: $d s = d' s \implies \text{row } s (P\text{-lower } d) = \text{row } s (P\text{-lower } d')$
<proof>

lemma *triangular-mat-P-upper*: *upper-triangular-mat* ($P\text{-upper } d$)
<proof>

lemma *slt-P-lower*: *strict-lower-triangular-mat* ($P\text{-lower } d$)
<proof>

lemma *lt-P-lower*: *lower-triangular-mat* ($P\text{-lower } d$)
<proof>

6.5.2 Gauss Seidel is a Regular Splitting

definition *Q-GS* $d = id\text{-blifun} - l *_R \mathcal{P}_L d$

definition *R-GS* $d = l *_R \mathcal{P}_U d$

lemma *splitting-gauss*: *is-splitting-blin* ($id\text{-blifun} - l *_R \mathcal{P}_1 (mk\text{-dec-det } d)$) (*Q-GS* d) (*R-GS* d)
<proof>

abbreviation *r-det_b* $d \equiv r\text{-dec}_b (mk\text{-dec-det } d)$

abbreviation *r-vec* $d \equiv \chi i. r\text{-dec}_b (mk\text{-dec-det } d) i$

abbreviation *Q-mat* $d \equiv blifun\text{-to-matrix} (Q\text{-GS } d)$

abbreviation *R-mat* $d \equiv blifun\text{-to-matrix} (R\text{-GS } d)$

lemma *Q-mat-def*: $Q\text{-mat } d = mat\ 1 - l *_R P\text{-lower } d$
<proof>

lemma *R-mat-def*: $R\text{-mat } d = l *_R P\text{-upper } d$
<proof>

lemma *triangular-mat-R*: *upper-triangular-mat* ($R\text{-mat } d$)
<proof>

definition $GS\text{-inv } d \ v \equiv \text{matrix-inv } (Q\text{-mat } d) *v (r\text{-vec } d + R\text{-mat } d *v \ v)$

$Q\text{-mat}$ can be expressed as an infinite sum of $P\text{-lower}$. It is therefore lower triangular.

lemma $\text{inv-}Q\text{-mat-suminf}$: $\text{matrix-inv } (Q\text{-mat } d) = (\sum k. (\text{matpow } (l *_R (P\text{-lower } d)) \ k))$
 $\langle \text{proof} \rangle$

lemma $\text{lt-}Q\text{-inv}$: $\text{lower-triangular-mat } (\text{matrix-inv } (Q\text{-mat } d))$
 $\langle \text{proof} \rangle$

Each row of the matrix $Q\text{-mat } d$ only depends on d 's actions in lower states.

lemma $\text{inv-}Q\text{-mat-indep}$:
assumes $\bigwedge i. i \leq s \implies d \ i = d' \ i \ i \leq s$
shows $\text{row } i \ (\text{matrix-inv } (Q\text{-mat } d)) = \text{row } i \ (\text{matrix-inv } (Q\text{-mat } d'))$
 $\langle \text{proof} \rangle$

As a result, also $GS\text{-inv}$ is independent of lower actions.

lemma $GS\text{-indep-high-states}$:
assumes $\bigwedge s'. s' \leq s \implies d \ s' = d' \ s'$
shows $GS\text{-inv } d \ v \ \$ \ s = GS\text{-inv } d' \ v \ \$ \ s$
 $\langle \text{proof} \rangle$

This recursive definition mimics the computation of the GS iteration.

lemma $GS\text{-inv-rec}$: $GS\text{-inv } d \ v = r\text{-vec } d + l *_R (P\text{-upper } d *v \ v + P\text{-lower } d *v \ (GS\text{-inv } d \ v))$
 $\langle \text{proof} \rangle$

lemma $\text{is-am-}GS\text{-inv-extend}$:
assumes $\bigwedge s. s < k \implies \text{is-arg-max } (\lambda d. GS\text{-inv } d \ v \ \$ \ s) (\lambda d. d \in D_D) \ d$
and $\text{is-arg-max } (\lambda a. GS\text{-inv } (d \ (k := a)) \ v \ \$ \ k) (\lambda a. a \in A \ k) \ a$
and $s \leq k$
and $d \in D_D$
shows $\text{is-arg-max } (\lambda d. GS\text{-inv } d \ v \ \$ \ s) (\lambda d. d \in D_D) \ (d \ (k := a))$
 $\langle \text{proof} \rangle$

lemma is-arg-max-GS-le :
 $\exists d. \forall s \leq k. \text{is-arg-max } (\lambda d. GS\text{-inv } d \ v \ \$ \ s) (\lambda d. d \in D_D) \ d$
 $\langle \text{proof} \rangle$

lemma ex-is-arg-max-GS :

$\exists d. \forall s. \text{is-arg-max } (\lambda d. \text{GS-inv } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
 $\langle \text{proof} \rangle$

function *GS-rec-fun* **where**

$\text{GS-rec-fun } v \ s = (\bigsqcup a \in A \ s. r \ (s, a) + l * ($
 $(\sum s' < s. \text{pmf } (K \ (s, a)) \ s' * (\text{GS-rec-fun } v \ s')) +$
 $(\sum s' \in \{s'. \ s \leq s'\}. \text{pmf } (K \ (s, a)) \ s' * v \ s'))$
 $\langle \text{proof} \rangle$

termination

$\langle \text{proof} \rangle$

declare *GS-rec-fun.simps*[*simp del*]

definition *GS-rec-elem* $v \ s \ a = r \ (s, a) + l * ($
 $(\sum s' < s. \text{pmf } (K \ (s, a)) \ s' * (\text{GS-rec-fun } v \ s')) +$
 $(\sum s' \in \{s'. \ s \leq s'\}. \text{pmf } (K \ (s, a)) \ s' * v \ s'))$

lemma *GS-rec-fun-elem*: $\text{GS-rec-fun } v \ s = (\bigsqcup a \in A \ s. \text{GS-rec-elem } v \ s \ a)$
 $\langle \text{proof} \rangle$

definition *GS-rec* $v = (\chi \ s. \text{GS-rec-fun } (\text{vec-nth } v) \ s)$

lemma *GS-rec-def'*: $\text{GS-rec } v \ \$ \ s = (\bigsqcup a \in A \ s. r \ (s, a) + l * ($
 $(\sum s' < s. \text{pmf } (K \ (s, a)) \ s' * (\text{GS-rec } v \ \$ \ s')) +$
 $(\sum s' \in \{s'. \ s \leq s'\}. \text{pmf } (K \ (s, a)) \ s' * v \ \$ \ s'))$
 $\langle \text{proof} \rangle$

lemma *GS-rec-eq*: $\text{GS-rec } v \ \$ \ s = (\bigsqcup a \in A \ s. r \ (s, a) + l * ($
 $(P\text{-lower } (d(s := a)) * v \ (\text{GS-rec } v)) \ \$ \ s + (P\text{-upper } (d(s := a)) * v$
 $v) \ \$ \ s))$
 $\langle \text{proof} \rangle$

definition *GS-rec-step* $d \ v \equiv r\text{-vec } d + l *_R \ (P\text{-lower } d * v \ \text{GS-rec } v$
 $+ P\text{-upper } d * v \ v)$

lemma *GS-rec-eq'*: $\text{GS-rec } v \ \$ \ s = (\bigsqcup a \in A \ s. \text{GS-rec-step } (d(s := a))$
 $v \ \$ \ s)$
 $\langle \text{proof} \rangle$

lemma *GS-rec-eq-vec*:

$\text{GS-rec } v \ \$ \ s = (\bigsqcup d \in D_D. \text{GS-rec-step } d \ v \ \$ \ s)$
 $\langle \text{proof} \rangle$

lift-definition *GS-rec-fun_b* $:: ('s \Rightarrow_b \text{real}) \Rightarrow ('s \Rightarrow_b \text{real}) \ \text{is } \text{GS-rec-fun}$
 $\langle \text{proof} \rangle$

definition *GS-rec-fun-inner* $(v :: 's \Rightarrow_b \text{real}) \ s \ a \equiv r \ (s, a) + l * ($
 $(\sum s' < s. \text{pmf } (K \ (s, a)) \ s' * (\text{GS-rec-fun}_b \ v \ s')) +$

$(\sum s' \in \{s'. s \leq s'\}. \text{pmf } (K (s,a)) s' * v s')$

definition *GS-rec-iter* **where**

$GS\text{-rec-iter } v s = (\bigsqcup a \in A s. r (s, a) + l * (\sum s' \in UNIV. \text{pmf } (K (s,a)) s' * v s'))$

lemma *GS-rec-fun-eq-GS-iter*:

assumes $\forall s' < s. v\text{-next } s' = GS\text{-rec-fun } v s' \forall s' \in \{s'. s \leq s'\}.$
 $v\text{-next } s' = v s'$

shows $GS\text{-rec-fun } v s = GS\text{-rec-iter } v\text{-next } s$

$\langle \text{proof} \rangle$

lemma *foldl-upd-notin*: $x \notin \text{set } X \implies \text{foldl } (\lambda f y. f(y := g f y)) c X$
 $x = c x$

$\langle \text{proof} \rangle$

lemma *foldl-upd-notin'*: $x \notin \text{set } Y \implies \text{foldl } (\lambda f y. f(y := g f y)) c$
 $(X @ Y) x = \text{foldl } (\lambda f y. f(y := g f y)) c X x$

$\langle \text{proof} \rangle$

lemma *sorted-list-of-set-split*:

assumes *finite* X

shows $\text{sorted-list-of-set } X = \text{sorted-list-of-set } \{x \in X. x < y\} @$
 $\text{sorted-list-of-set } \{x \in X. y \leq x\}$

$\langle \text{proof} \rangle$

lemma *sorted-list-of-set-split'*:

assumes *finite* X

shows $\text{sorted-list-of-set } X = \text{sorted-list-of-set } \{x \in X. x \leq y\} @$
 $\text{sorted-list-of-set } \{x \in X. y < x\}$

$\langle \text{proof} \rangle$

lemma *GS-rec-fun-code*: $GS\text{-rec-fun } v s = \text{foldl } (\lambda v s. v(s := GS\text{-rec-iter}$
 $v s)) v (\text{sorted-list-of-set } \{..s\}) s$

$\langle \text{proof} \rangle$

lemma *GS-rec-fun-code'*: $GS\text{-rec-fun } v s = \text{foldl } (\lambda v s. v(s := GS\text{-rec-iter}$
 $v s)) v (\text{sorted-list-of-set } UNIV) s$

$\langle \text{proof} \rangle$

lemma *GS-rec-fun-code''*: $GS\text{-rec-fun } v = \text{foldl } (\lambda v s. v(s := GS\text{-rec-iter}$
 $v s)) v (\text{sorted-list-of-set } UNIV)$

$\langle \text{proof} \rangle$

lemma *GS-rec-eq-elem*: $GS\text{-rec } v \$ s = GS\text{-rec-fun } (vec\text{-nth } v) s$

$\langle \text{proof} \rangle$

lemma *GS-rec-step-elem*: $GS\text{-rec-step } d \ v \ \$ \ s = r \ (s, d \ s) + l * ((\sum s' < s. \text{pmf } (K \ (s, d \ s)) \ s' * GS\text{-rec } v \ \$ \ s') + (\sum s' \in \{s'. \ s \leq s'\}. \text{pmf } (K \ (s, d \ s)) \ s' * v \ \$ \ s'))$
 ⟨proof⟩

lemma *is-arg-max-GS-rec-step-act*:
assumes $d \in D_D$ *is-arg-max* $(\lambda a. GS\text{-rec-step } (d'(s := a)) \ v \ \$ \ s)$ $(\lambda a. a \in A \ s) \ a$
shows *is-arg-max* $(\lambda d. GS\text{-rec-step } d \ v \ \$ \ s)$ $(\lambda d. d \in D_D) \ (d(s := a))$
 ⟨proof⟩

lemma *is-arg-max-GS-rec-step-act'*:
assumes $d \in D_D$ *is-arg-max* $(\lambda a. GS\text{-rec-step } (d'(s := a)) \ v \ \$ \ s)$ $(\lambda a. a \in A \ s) \ (d \ s)$
shows *is-arg-max* $(\lambda d. GS\text{-rec-step } d \ v \ \$ \ s)$ $(\lambda d. d \in D_D) \ d$
 ⟨proof⟩

lemma
is-arg-max-GS-rec:
assumes $\bigwedge s. \text{is-arg-max } (\lambda d. GS\text{-rec-step } d \ v \ \$ \ s)$ $(\lambda d. d \in D_D) \ d$
shows $GS\text{-rec } v = GS\text{-rec-step } d \ v$
 ⟨proof⟩

lemma
is-arg-max-GS-rec':
assumes *is-arg-max* $(\lambda d. GS\text{-rec-step } d \ v \ \$ \ s)$ $(\lambda d. d \in D_D) \ d$
shows $GS\text{-rec } v \ \$ \ s = GS\text{-rec-step } d \ v \ \$ \ s$
 ⟨proof⟩

lemma
GS-rec-eq-GS-inv:
assumes $\bigwedge s. \text{is-arg-max } (\lambda d. GS\text{-rec-step } d \ v \ \$ \ s)$ $(\lambda d. d \in D_D) \ d$
shows $GS\text{-rec } v = GS\text{-inv } d \ v$
 ⟨proof⟩

lemma
GS-rec-step-eq-GS-inv:
assumes $\bigwedge s. \text{is-arg-max } (\lambda d. GS\text{-rec-step } d \ v \ \$ \ s)$ $(\lambda d. d \in D_D) \ d$
shows $GS\text{-rec-step } d \ v = GS\text{-inv } d \ v$
 ⟨proof⟩

lemma *strict-lower-triangular-mat-mult*:
assumes *strict-lower-triangular-mat* $M \ \bigwedge i. \ i < j \implies v \ \$ \ i = v' \ \$ \ i$
shows $(M * v \ v) \ \$ \ j = (M * v \ v') \ \$ \ j$
 ⟨proof⟩

lemma *Q-mat-invertible*: *invertible* $(Q\text{-mat } d)$
 ⟨proof⟩

lemma *GS-eq-GS-inv*:

assumes $\bigwedge s. s \leq k \implies \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
assumes $s \leq k$
shows $\text{GS-rec-step } d \ v \ \$ \ s = \text{GS-inv } d \ v \ \$ \ s$
<proof>

lemma *is-arg-max-GS-imp-splitting'*:

assumes $\bigwedge s. s \leq k \implies \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
assumes $s \leq k$
shows $\text{is-arg-max } (\lambda d. \text{GS-inv } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
<proof>

lemma *is-am-GS-rec-step-indep*:

assumes $d \ s = d' \ s$
assumes $\text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
shows $\text{GS-rec } v \ \$ \ s = \text{GS-rec-step } d' \ v \ \$ \ s$
<proof>

lemma *is-am-GS-rec-step-indep'*:

assumes $d \ s = d' \ s$
assumes $\text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
shows $\text{GS-rec } v \ \$ \ s = \text{GS-rec-step } d' \ v \ \$ \ s$
<proof>

lemma *is-arg-max-GS-imp-splitting''*:

assumes $\bigwedge s. s \leq k \implies \text{is-arg-max } (\lambda d. \text{GS-inv } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
assumes $s \leq k$
shows $\text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d \wedge \text{GS-inv } d \ v \ \$ \ s = \text{GS-rec } v \ \$ \ s$
<proof>

lemma *is-arg-max-GS-imp-splitting'''*:

assumes $\bigwedge s. s \leq k \implies \text{is-arg-max } (\lambda d. \text{GS-inv } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
assumes $s \leq k$
shows $\text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
<proof>

lemma *is-arg-max-GS-imp-splitting*:

assumes $\bigwedge s. \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
shows $\text{is-arg-max } (\lambda d. \text{GS-inv } d \ v \ \$ \ k) \ (\lambda d. d \in D_D) \ d$
<proof>

lemma *is-arg-max-gs-iff*:

assumes $d \in D_D$

shows $(\forall s \leq k. \text{is-arg-max } (\lambda d. \text{GS-inv } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d)$
 \longleftrightarrow
 $(\forall s \leq k. \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d)$
 $\langle \text{proof} \rangle$

lemma *GS-opt-indep-high*:

assumes $(\bigwedge s'. s' < s \implies \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s') \text{ is-dec-det } d) \ s' < s \ a \in A \ s$
shows $\text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s') \ \text{is-dec-det } (d(s := a))$
 $\langle \text{proof} \rangle$

lemma *mult-mat-vec-nth*: $(X * v \ x) \ \$ \ i = \text{scalar-product } (\text{row } i \ X) \ x$
 $\langle \text{proof} \rangle$

lemma *ext-GS-opt-le*:

assumes $(\bigwedge s'. s' < s \implies \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s') \ (\lambda d. d \in D_D) \ d)$
and $\text{is-arg-max } (\lambda a. \text{GS-rec-step } (d(s := a)) \ v \ \$ \ s) \ (\lambda a. a \in A \ s)$
 $a \ s' \leq s$
and $d \in D_D$
shows $\text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s') \ (\lambda d. d \in D_D) \ (d(s := a))$
 $\langle \text{proof} \rangle$

lemma *ex-GS-opt-le*:

shows $\exists d. (\forall s' \leq s. \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s') \ (\lambda d. d \in D_D) \ d)$
 $\langle \text{proof} \rangle$

lemma *ex-GS-opt*:

shows $\exists d. \forall s. \text{is-arg-max } (\lambda d. \text{GS-rec-step } d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d$
 $\langle \text{proof} \rangle$

lemma *GS-rec-eq-GS-inv'*: $\text{GS-rec } v \ \$ \ s = (\bigsqcup d \in D_D. \text{GS-inv } d \ v \ \$ \ s)$
 $\langle \text{proof} \rangle$

lemma *GS-rec-fun-eq-GS-inv*: $\text{GS-rec-fun } v \ s = (\bigsqcup d \in D_D. \text{GS-inv } d \ (\text{vec-lambda } v) \ \$ \ s)$
 $\langle \text{proof} \rangle$

lemma *invertible-Q-GS*: $\text{invertible}_L \ (Q\text{-GS } d)$ **for** d
 $\langle \text{proof} \rangle$

lemma *ex-opt-blifun*: $\exists d. \forall s. \text{is-arg-max } (\lambda d. ((\text{inv}_L \ (Q\text{-GS } d)) \ (r\text{-det}_b \ d + (R\text{-GS } d) \ v)) \ s) \ \text{is-dec-det } d$
 $\langle \text{proof} \rangle$

lemma *GS-inv-blinfun-to-matrix*: $((inv_L (Q-GS d)) (r-det_b d + R-GS d v)) = Bfun (vec-nth (GS-inv d (vec-lambda v)))$
 ⟨proof⟩

lemma *norm-GS-QR-le-disc*: $norm (inv_L (Q-GS d) o_L R-GS d) \leq l$
 ⟨proof⟩

sublocale *GS: MDP-QR A K r l Q-GS R-GS*
rewrites $GS.\mathcal{L}_b-split = GS-rec-fun_b$
 ⟨proof⟩

abbreviation *gs-measure* $\equiv (\lambda(eps, v).$
 if $v = \nu_b-opt \vee l = 0$
 then 0
 else $nat (ceiling (log (1/l) (dist v \nu_b-opt) - log (1/l) (eps * (1-l)) / (8 * l))))$

lemma *dist-L_b-split-lt-dist-opt*: $dist v (GS-rec-fun_b v) \leq 2 * dist v \nu_b-opt$
 ⟨proof⟩

lemma *GS-QR-disc-le-disc*: $GS.QR-disc \leq l$
 ⟨proof⟩

lemma *gs-rel-dec*:
assumes $l \neq 0$ $GS-rec-fun_b v \neq \nu_b-opt$
shows $\lceil log (1 / l) (dist (GS-rec-fun_b v) \nu_b-opt) - c \rceil < \lceil log (1 / l) (dist v \nu_b-opt) - c \rceil$
 ⟨proof⟩

function *gs-iteration* :: $real \Rightarrow ('s \Rightarrow_b real) \Rightarrow ('s \Rightarrow_b real)$ **where**
 gs-iteration $eps v =$
 (if $2 * l * dist v (GS-rec-fun_b v) < eps * (1-l) \vee eps \leq 0$ *then*
 $GS-rec-fun_b v$ *else* $gs-iteration eps (GS-rec-fun_b v)$)
 ⟨proof⟩

termination
 ⟨proof⟩

lemma *THE-fix-GS-rec-fun_b*: $(THE v. GS-rec-fun_b v = v) = \nu_b-opt$
 ⟨proof⟩

The distance between an estimate for the value and the optimal value can be bounded with respect to the distance between the estimate and the result of applying it to \mathcal{L}_b

lemma *contraction-L-split-dist*: $(1 - l) * dist v \nu_b-opt \leq dist v (GS-rec-fun_b v)$
 ⟨proof⟩

lemma *dist- \mathcal{L}_b -split-opt-eps*:

assumes $eps > 0 \ 2 * l * dist \ v \ (GS-rec-fun_b \ v) < eps * (1-l)$

shows $dist \ (GS-rec-fun_b \ v) \ \nu_b-opt < eps / 2$

<proof>

end

context *MDP-ord*

begin

lemma *is-am-GS-inv-extend'*:

assumes $(\bigwedge s. s < x \implies is-arg-max \ (\lambda d. GS-inv \ d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ d)$

assumes $is-arg-max \ (\lambda d. GS-rec-step \ d \ v \ \$ \ x) \ (\lambda d. d \in D_D) \ (d(x := a))$

assumes $s \leq x \ d \in D_D$

shows $is-arg-max \ (\lambda d. GS-inv \ d \ v \ \$ \ s) \ (\lambda d. d \in D_D) \ (d(x := a))$

<proof>

definition *opt-policy-gs'* $d \ v \ s = (LEAST \ a. is-arg-max \ (\lambda a. GS-rec-step \ (d(s := a)) \ v \ \$ \ s) \ (\lambda a. a \in A \ s) \ a)$

definition *GS-iter* $a \ v \ s = r \ (s, a) + l * (\sum s' \in UNIV. pmf \ (K(s, a)) \ s' * v \ \$ \ s')$

definition *GS-iter-max* $v \ s = (\bigsqcup a \in A \ s. GS-iter \ a \ v \ s)$

lemma *GS-rec-eq-iter*:

assumes $\bigwedge s. s < k \implies v' \ \$ \ s = GS-rec \ v \ \$ \ s \ \bigwedge s. k \leq s \implies v' \ \$ \ s = v \ \$ \ s$

shows $GS-rec-step \ (d(k := a)) \ v \ \$ \ k = GS-iter \ a \ v' \ k$

<proof>

lemma *GS-rec-eq-iter-max*:

assumes $\bigwedge s. s < k \implies v' \ \$ \ s = GS-rec \ v \ \$ \ s \ \bigwedge s. k \leq s \implies v' \ \$ \ s = v \ \$ \ s$

shows $GS-rec \ v \ \$ \ k = GS-iter-max \ v' \ k$

<proof>

definition *GS-iter-arg-max* $v \ s = (LEAST \ a. is-arg-max \ (\lambda a. GS-iter \ a \ v \ s) \ (\lambda a. a \in A \ s) \ a)$

definition *GS-rec-am-code* $v \ d \ s = foldl \ (\lambda v d s. vd(s := (GS-iter-max \ (\chi \ s. fst \ (vd \ s)) \ s, \ GS-iter-arg-max \ (\chi \ s. fst \ (vd \ s)) \ s))) \ (\lambda s. (v \ \$ \ s, d \ s)) \ (sorted-list-of-set \ \{..s\}) \ s$

definition *GS-rec-am-code'* $v \ d \ s = foldl \ (\lambda v d s. vd(s := (GS-iter-max \ (\chi \ s. fst \ (vd \ s)) \ s, \ GS-iter-arg-max \ (\chi \ s. fst \ (vd \ s)) \ s))) \ (\lambda s. (v \ \$ \ s, d \ s)) \ (sorted-list-of-set \ UNIV) \ s$

lemma *GS-rec-am-code'*: $GS\text{-rec-am-code} = GS\text{-rec-am-code}'$
 ⟨proof⟩

lemma *opt-policy-gs'-eq-GS-iter*:
assumes $\bigwedge s. s < k \implies v' \$ s = GS\text{-rec } v \$ s \wedge s. k \leq s \implies v' \$ s = v \$ s$
shows $opt\text{-policy-gs}' d v k = GS\text{-iter-arg-max } v' k$
 ⟨proof⟩

lemma *opt-policy-gs'-eq-GS-iter'*:
 $opt\text{-policy-gs}' d v k = GS\text{-iter-arg-max } (\chi s. \text{if } s < k \text{ then } GS\text{-rec } v \$ s \text{ else } v \$ s) k$
 ⟨proof⟩

lemma *opt-policy-gs'-is-dec-det*: $opt\text{-policy-gs}' d v \in D_D$
 ⟨proof⟩

lemma *opt-policy-gs'-is-arg-max*: $is\text{-arg-max } (\lambda d. GS\text{-inv } d v \$ s) (\lambda d. d \in D_D) (opt\text{-policy-gs}' d v)$
 ⟨proof⟩

lemma *GS-rec-am-code* $v d s = (GS\text{-rec } v \$ s, opt\text{-policy-gs}' d v s)$
 ⟨proof⟩

lemma *GS-rec-am-code-eq*: $GS\text{-rec-am-code } v d s = (GS\text{-rec } v \$ s, opt\text{-policy-gs}' d v s)$
 ⟨proof⟩

definition *GS-rec-iter-arg-max* **where**

$GS\text{-rec-iter-arg-max } v s = (LEAST a. is\text{-arg-max } (\lambda a. r (s, a) + l * (\sum s' \in UNIV. pmf (K (s, a)) s' * v s')) (\lambda a. a \in A s) a)$

definition *opt-policy-gs* $v s = (LEAST a. is\text{-arg-max } (\lambda a. GS\text{-rec-fun-inner } v s a) (\lambda a. a \in A s) a)$

lemma *opt-policy-gs-eq'*: $opt\text{-policy-gs } v = opt\text{-policy-gs}' d (vec\text{-lambda } v)$
 ⟨proof⟩

declare *gs-iteration.simps*[simp del]

lemma *gs-iteration-error*:
assumes $eps > 0$
shows $dist (gs\text{-iteration } eps v) \nu_b\text{-opt} < eps / 2$
 ⟨proof⟩

lemma *GS-rec-fun-inner-vec*: $GS\text{-rec-fun-inner } v s a = GS\text{-rec-step } (d(s := a)) (vec\text{-lambda } v) \$ s$
 ⟨proof⟩

lemma *find-policy-error-bound-gs*:
assumes $\text{eps} > 0 \ 2 * l * \text{dist } v \ (GS\text{-rec-fun}_b \ v) < \text{eps} * (1-l)$
shows $\text{dist } (\nu_b \ (\text{mk-stationary-det } (\text{opt-policy-gs } (GS\text{-rec-fun}_b \ v))))$
 $\nu_b\text{-opt} < \text{eps}$
 $\langle \text{proof} \rangle$

definition *vi-gs-policy* $\text{eps } v = \text{opt-policy-gs } (\text{gs-iteration } \text{eps } v)$

lemma *vi-gs-policy-opt*:
assumes $0 < \text{eps}$
shows $\text{dist } (\nu_b \ (\text{mk-stationary-det } (\text{vi-gs-policy } \text{eps } v))) \nu_b\text{-opt} < \text{eps}$
 $\langle \text{proof} \rangle$

lemma *GS-rec-iter-eq-iter-max*: $GS\text{-rec-iter } v = GS\text{-iter-max } (\text{vec-lambda}$
 $v)$
 $\langle \text{proof} \rangle$
end

end

theory *Algorithms*
imports
Value-Iteration
Policy-Iteration
Modified-Policy-Iteration
Splitting-Methods
begin
end

theory *Code-DP*
imports
Value-Iteration
Policy-Iteration
Modified-Policy-Iteration
Splitting-Methods

HOL-Library.Code-Target-Numeral
Gauss-Jordan.Code-Generation-IArrays
begin

7 Code Generation for MDP Algorithms

7.1 Least Argmax

lemma *least-list*:
assumes $\text{sorted } xs \ \exists x \in \text{set } xs. P \ x$

shows $(LEAST\ x \in\ set\ xs.\ P\ x) = the\ (find\ P\ xs)$
 $\langle proof \rangle$

definition $least-enum\ P = the\ (find\ P\ (sorted-list-of-set\ (UNIV :: ('b::\ {finite,\ linorder})\ set)))$

lemma $least-enum-eq:\ \exists\ x.\ P\ x \implies least-enum\ P = (LEAST\ x.\ P\ x)$
 $\langle proof \rangle$

definition $least-max-arg-max-list\ f\ init\ xs = foldl\ (\lambda(am,\ m)\ x.\ if\ f\ x > m\ then\ (x,\ f\ x)\ else\ (am,\ m))\ init\ xs$

lemma $snd-least-max-arg-max-list:$
 $snd\ (least-max-arg-max-list\ f\ (n,\ f\ n)\ xs) = (MAX\ x \in\ insert\ n\ (set\ xs).\ f\ x)$
 $\langle proof \rangle$

lemma $least-max-arg-max-list-snd-fst:\ snd\ (least-max-arg-max-list\ f\ (x,\ f\ x)\ xs) = f\ (fst\ (least-max-arg-max-list\ f\ (x,\ f\ x)\ xs))$
 $\langle proof \rangle$

lemma $fst-least-max-arg-max-list:$
fixes $f :: - \Rightarrow - :: linorder$
assumes $sorted\ (n\#\ xs)$
shows $fst\ (least-max-arg-max-list\ f\ (n,\ f\ n)\ xs) = (LEAST\ x.\ is-arg-max\ f\ (\lambda x.\ x \in\ insert\ n\ (set\ xs))\ x)$
 $\langle proof \rangle$

definition $least-arg-max-enum\ f\ X = (let\ xs = sorted-list-of-set\ (X :: (- ::\ {finite,\ linorder})\ set)\ in\ fst\ (least-max-arg-max-list\ f\ (hd\ xs,\ f\ (hd\ xs))\ (tl\ xs)))$

definition $least-max-arg-max-enum\ f\ X = (let\ xs = sorted-list-of-set\ (X :: (- ::\ {finite,\ linorder})\ set)\ in\ (least-max-arg-max-list\ f\ (hd\ xs,\ f\ (hd\ xs))\ (tl\ xs)))$

lemma $least-arg-max-enum-correct:$
assumes $X \neq \{\}$
shows
 $(least-arg-max-enum\ (f :: - \Rightarrow (- :: linorder))\ X) = (LEAST\ x.\ is-arg-max\ f\ (\lambda x.\ x \in\ X)\ x)$
 $\langle proof \rangle$

lemma $least-max-arg-max-enum-correct1:$
assumes $X \neq \{\}$
shows $fst\ (least-max-arg-max-enum\ (f :: - \Rightarrow (- :: linorder))\ X) = (LEAST\ x.\ is-arg-max\ f\ (\lambda x.\ x \in\ X)\ x)$
 $\langle proof \rangle$

lemma *least-max-arg-max-enum-correct2*:
assumes $X \neq \{\}$
shows $\text{snd } (\text{least-max-arg-max-enum } (f :: - \Rightarrow (- :: \text{linorder})) X) =$
 $(\text{MAX } x \in X. f x)$
 $\langle \text{proof} \rangle$

7.2 Functions as Vectors

typedef $('a, 'b) \text{ Fun} = \text{UNIV} :: ('a \Rightarrow 'b) \text{ set}$
 $\langle \text{proof} \rangle$

setup-lifting *type-definition-Fun*

lift-definition $\text{to-Fun} :: ('a \Rightarrow 'b) \Rightarrow ('a, 'b) \text{ Fun}$ **is** id $\langle \text{proof} \rangle$

definition $\text{fun-to-vec } (v :: ('a :: \text{finite}, 'b) \text{ Fun}) = \text{vec-lambda } (\text{Rep-Fun } v)$

lift-definition $\text{vec-to-fun} :: 'b \wedge 'a \Rightarrow ('a, 'b) \text{ Fun}$ **is** vec-nth $\langle \text{proof} \rangle$

lemma $\text{Fun-inverse}[\text{simp}] : \text{Rep-Fun } (\text{Abs-Fun } f) = f$
 $\langle \text{proof} \rangle$

lift-definition $\text{zero-Fun} :: ('a, 'b :: \text{zero}) \text{ Fun}$ **is** 0 $\langle \text{proof} \rangle$

code-datatype *vec-to-fun*

lemmas *vec-to-fun.rep-eq* $[\text{code}]$

instantiation $\text{Fun} :: (\text{enum}, \text{equal}) \text{ equal}$

begin

definition $\text{equal-Fun } (f :: ('a :: \text{enum}, 'b :: \text{equal}) \text{ Fun}) g = (\text{Rep-Fun } f$
 $= \text{Rep-Fun } g)$

instance

$\langle \text{proof} \rangle$

end

7.3 Bounded Functions as Vectors

lemma $\text{Bfun-inverse-fin}[\text{simp}] : \text{apply-bfun } (\text{Bfun } (f :: 'c :: \text{finite} \Rightarrow -))$
 $= f$
 $\langle \text{proof} \rangle$

definition $\text{bfun-to-vec } (v :: ('a :: \text{finite}) \Rightarrow_b ('b :: \text{metric-space})) = \text{vec-lambda}$
 v

definition $\text{vec-to-bfun } v = \text{Bfun } (\text{vec-nth } v)$

code-datatype *vec-to-bfun*

lemma *apply-bfun-vec-to-bfun*[code]: *apply-bfun (vec-to-bfun f) x = f*
 $\$ x$
 $\langle proof \rangle$

lemma [code]: $0 = \text{vec-to-bfun } 0$
 $\langle proof \rangle$

7.4 IArrays with Lengths in the Type

typedef ($'s :: \text{mod-type}, 'a$) *iarray-type* = $\{\text{arr} :: 'a \text{ iarray}. \text{IArray.length arr} = \text{CARD}('s)\}$
 $\langle proof \rangle$

setup-lifting *type-definition-iarray-type*

lift-definition *fun-to-iarray-t* :: ($'s :: \{\text{mod-type}\} \Rightarrow 'a$) $\Rightarrow ('s, 'a) \text{ iarray-type}$ **is** $\lambda f. \text{IArray.of-fun } (\lambda s. f (\text{from-nat } s)) (\text{CARD}('s))$
 $\langle proof \rangle$

lift-definition *iarray-t-sub* :: ($'s :: \text{mod-type}, 'a$) *iarray-type* $\Rightarrow 's \Rightarrow 'a$ **is** $\lambda v x. \text{IArray.sub } v (\text{to-nat } x) \langle proof \rangle$

lift-definition *iarray-to-vec* :: ($'s, 'a$) *iarray-type* $\Rightarrow 'a \wedge 's :: \{\text{mod-type}, \text{finite}\}$ **is** $\lambda v. (\chi s. \text{IArray.sub } v (\text{to-nat } s)) \langle proof \rangle$

lift-definition *vec-to-iarray* :: $'a \wedge 's :: \{\text{mod-type}, \text{finite}\} \Rightarrow ('s, 'a) \text{ iarray-type}$ **is** $\lambda v. \text{IArray.of-fun } (\lambda s. v \$ ((\text{from-nat } s) :: 's)) (\text{CARD}('s))$
 $\langle proof \rangle$

lemma *length-iarray-type* [simp]: $\text{length } (\text{IArray.list-of } (\text{Rep-iarray-type } (v :: ('s :: \{\text{mod-type}\}, 'a) \text{ iarray-type}))) = \text{CARD}('s)$
 $\langle proof \rangle$

lemma *iarray-t-eq-iff*: $(v = w) = (\forall x. \text{iarray-t-sub } v x = \text{iarray-t-sub } w x)$
 $\langle proof \rangle$

lemma *iarray-to-vec-inv*: $\text{iarray-to-vec } (\text{vec-to-iarray } v) = v$
 $\langle proof \rangle$

lemma *vec-to-iarray-inv*: $\text{vec-to-iarray } (\text{iarray-to-vec } v) = v$
 $\langle proof \rangle$

code-datatype *iarray-to-vec*

lemma *vec-nth-iarray-to-vec*[code]: $\text{vec-nth } (\text{iarray-to-vec } v) x = \text{iarray-t-sub } v x$

$\langle \text{proof} \rangle$

lemma *vec-lambda-iarray-t*[code]: $\text{vec-lambda } v = \text{iarray-to-vec } (\text{fun-to-iarray-t } v)$
 $\langle \text{proof} \rangle$

lemma *zero-iarray*[code]: $0 = \text{iarray-to-vec } (\text{fun-to-iarray-t } 0)$
 $\langle \text{proof} \rangle$

7.5 Value Iteration

locale *vi-code* =
 $\text{MDP-ord } A \ K \ r \ l$ **for** $A :: 's::\text{mod-type} \Rightarrow ('a::\{\text{finite, wellorder}\})$
 set
 and $K :: ('s::\{\text{finite, mod-type}\} \times 'a::\{\text{finite, wellorder}\}) \Rightarrow 's \ \text{pmf}$
 and $r \ l$
begin
definition *vi-test* ($v::'s \Rightarrow_b \ \text{real}$) $v' \ \text{eps} = 2 * l * \text{dist } v \ v'$

partial-function (*tailrec*) *value-iteration-partial* **where** [code]: *value-iteration-partial*
 $\text{eps } v =$
 ($\text{let } v' = \mathcal{L}_b \ v \ \text{in}$
 ($\text{if } 2 * l * \text{dist } v \ v' < \text{eps} * (1 - l)$ then v' else ($\text{value-iteration-partial}$
 $\text{eps } v'$)))

lemma *vi-eq-partial*: $\text{eps} > 0 \implies \text{value-iteration-partial } \text{eps } v =$
 $\text{value-iteration } \text{eps } v$
 $\langle \text{proof} \rangle$

definition *L-det* $d = L \ (\text{mk-dec-det } d)$

lemma *code-L-det* [code]: $L\text{-det } d \ (\text{vec-to-bfun } v) = \text{vec-to-bfun } (\chi \ s.$
 $L_a \ (d \ s) \ (\text{vec-nth } v) \ s)$
 $\langle \text{proof} \rangle$

lemma *code-L_b* [code]: $\mathcal{L}_b \ (\text{vec-to-bfun } v) = \text{vec-to-bfun } (\chi \ s. (\text{MAX } a$
 $\in A \ s. \ r \ (s, a) + l * \text{measure-pmf.expectation } (K \ (s, a)) \ (\text{vec-nth } v)))$
 $\langle \text{proof} \rangle$

lemma *code-value-iteration*[code]: $\text{value-iteration } \text{eps} \ (\text{vec-to-bfun } v)$
 $=$
 ($\text{if } \text{eps} \leq 0$ then $\mathcal{L}_b \ (\text{vec-to-bfun } v)$ else $\text{value-iteration-partial } \text{eps}$
 $(\text{vec-to-bfun } v)$)
 $\langle \text{proof} \rangle$

lift-definition *find-policy-impl* :: $('s \Rightarrow_b \ \text{real}) \Rightarrow ('s, 'a) \ \text{Fun} \ \text{is } \lambda v.$
 $\text{find-policy}' \ v \langle \text{proof} \rangle$

lemma *code-find-policy-impl*: $\text{find-policy-impl } v = \text{vec-to-fun } (\chi \ s.$
 $(\text{LEAST } x. \ x \in \text{opt-acts } v \ s))$

$\langle \text{proof} \rangle$

lemma *code-find-policy-impl-opt*[code]: *find-policy-impl* $v = \text{vec-to-fun}$
(χ s . *least-arg-max-enum* (λa . L_a a v s) (A s))
 $\langle \text{proof} \rangle$

lemma *code-vi-policy'*[code]: *vi-policy'* eps $v = \text{Rep-Fun}$ (*find-policy-impl*
(*value-iteration* eps v))
 $\langle \text{proof} \rangle$

7.6 Policy Iteration

partial-function (*tailrec*) *policy-iteration-partial* **where** [code]: *policy-iteration-partial* $d =$
(*let* $d' = \text{policy-step } d$ *in if* $d = d'$ *then* d *else* *policy-iteration-partial*
 d')

lemma *pi-eq-partial*: $d \in D_D \implies \text{policy-iteration-partial } d = \text{policy-iteration } d$
 $\langle \text{proof} \rangle$

definition *P-mat* $d = (\chi$ i j . *pmf* (K (i , *Rep-Fun* d i)) j)

definition *r-vec'* $d = (\chi$ i . $r(i, \text{Rep-Fun } d$ i))

lift-definition *policy-eval'* :: ($'s :: \{\text{mod-type}, \text{finite}\}$, $'a$) *Fun* \Rightarrow ($'s \Rightarrow_b$
real) **is** *policy-eval* $\langle \text{proof} \rangle$

lemma *mat-eq-blinfun*: *mat* $1 - l *_R$ (*P-mat* (*vec-to-fun* d)) = *blinfun-to-matrix*
(*id-blinfun* $- l *_R \mathcal{P}_1$ (*mk-dec-det* (*vec-nth* d))))
 $\langle \text{proof} \rangle$

lemma ν_b -*vec*: *policy-eval'* (*vec-to-fun* d) = *vec-to-bfun* (*matrix-inv*
(*mat* $1 - l *_R$ (*P-mat* (*vec-to-fun* d))) $*v$ (*r-vec'* (*vec-to-fun* d))))
 $\langle \text{proof} \rangle$

lemma ν_b -*vec-opt*[code]: *policy-eval'* (*vec-to-fun* d) = *vec-to-bfun* (*Matrix-To-IArray.iarray-to-vec*
(*Matrix-To-IArray.vec-to-iarray* ((*fst* (*Gauss-Jordan-PA* ((*mat* $1 - l *_R$ (*P-mat* (*vec-to-fun* d)))))) $*v$ (*r-vec'* (*vec-to-fun* d))))))
 $\langle \text{proof} \rangle$

lift-definition *policy-improvement'* :: ($'s$, $'a$) *Fun* \Rightarrow ($'s \Rightarrow_b$ *real*) \Rightarrow
($'s$, $'a$) *Fun*
is *policy-improvement* $\langle \text{proof} \rangle$

lemma [code]: *policy-improvement'* (*vec-to-fun* d) $v = \text{vec-to-fun}$ (χ
 s . (*if is-arg-max* (λa . L_a a v s) (λa . $a \in A$ s) (d $\$$ s) *then* d $\$$ s *else*
LEAST x . *is-arg-max* (λa . L_a a v s) (λa . $a \in A$ s) x))
 $\langle \text{proof} \rangle$

lift-definition *policy-step'* :: ('s, 'a) Fun \Rightarrow ('s, 'a) Fun
is *policy-step* \langle proof \rangle

lemma [code]: *policy-step' d = policy-improvement' d (policy-eval' d)*
 \langle proof \rangle

lift-definition *policy-iteration-partial'* :: ('s, 'a) Fun \Rightarrow ('s, 'a) Fun
is *policy-iteration-partial* \langle proof \rangle

lemma [code]: *policy-iteration-partial' d = (let d' = policy-step' d in if d = d' then d else policy-iteration-partial' d')*
 \langle proof \rangle

lift-definition *policy-iteration'* :: ('s, 'a) Fun \Rightarrow ('s, 'a) Fun **is** *policy-iteration* \langle proof \rangle

lemma *code-policy-iteration'*[code]: *policy-iteration' d = (if Rep-Fun d \notin D_D then d else (policy-iteration-partial' d))*
 \langle proof \rangle

lemma *code-policy-iteration*[code]: *policy-iteration d = Rep-Fun (policy-iteration' (vec-to-fun (vec-lambda d)))*
 \langle proof \rangle

7.7 Gauss-Seidel Iteration

partial-function (*tailrec*) *gs-iteration-partial* **where**
 [code]: *gs-iteration-partial eps v = (let v' = (GS-rec-fun_b v) in (if 2 * l * dist v v' < eps * (1 - l) then v' else gs-iteration-partial eps v'))*

lemma *gs-iteration-partial-eq*: *eps > 0 \implies gs-iteration-partial eps v = gs-iteration eps v*
 \langle proof \rangle

lemma *gs-iteration-code-opt*[code]: *gs-iteration eps v = (if eps \leq 0 then GS-rec-fun_b v else gs-iteration-partial eps v)*
 \langle proof \rangle

definition *vec-upd v i x = (χ j. if i = j then x else v \$ j)*

lemma *GS-rec-eq-fold*: *GS-rec v = foldl ($\lambda v s. (vec-upd v s (GS-iter-max v s))) v$ (sorted-list-of-set UNIV)*
 \langle proof \rangle

lemma *GS-rec-fun-code''''*[code]: *GS-rec-fun_b (vec-to-bfun v) = vec-to-bfun (foldl ($\lambda v s. (vec-upd v s (GS-iter-max v s))) v$ (sorted-list-of-set*

UNIV)
 ⟨proof⟩

lemma *GS-iter-max-code* [code]: *GS-iter-max* v s = (*MAX* $a \in A$ s .
GS-iter a v s)
 ⟨proof⟩

lift-definition *opt-policy-gs''* :: ($'s \Rightarrow_b \text{real}$) \Rightarrow ($'s, 'a$) *Fun is opt-policy-gs*⟨proof⟩

declare *opt-policy-gs''.rep-eq*[*symmetric, code*]

lemma *GS-rec-am-code'-prod*: *GS-rec-am-code'* v d =
 ($\lambda s'$. (
 let (v', d') = *foldl* ($\lambda(v,d) s$. ($v(s := (\text{GS-iter-max } (\text{vec-lambda } v) s)$), $d(s := \text{GS-iter-arg-max } (\text{vec-lambda } v) s)$)) (*vec-nth* v , d)
 (*sorted-list-of-set UNIV*)
 in ($v' s', d' s'$)))
 ⟨proof⟩

lemma *code-GS-rec-am-arr-opt*[code]: *opt-policy-gs''* (*vec-to-bfun* v) =
vec-to-fun ((*snd* (*foldl* ($\lambda(v, d) s$.
 let (am, m) = *least-max-arg-max-enum* ($\lambda a. r(s, a) + l * (\sum s' \in$
 UNIV. pmf ($K(s, a) s' * v \$ s'$)) ($A s$) *in*
 (*vec-upd* $v s m$, *vec-upd* $d s am$))
 ($v, (\chi s. (\text{least-enum } (\lambda a. a \in A s)))$) (*sorted-list-of-set UNIV*))))
 ⟨proof⟩

7.8 Modified Policy Iteration

sublocale *MDP-MPI* $A K r l \lambda X$. *Least* ($\lambda x. x \in X$)
 ⟨proof⟩

definition $d0$ s = (*LEAST* $a. a \in A s$)

lift-definition $d0'$:: ($'s, 'a$) *Fun is d0*⟨proof⟩

lemma *d0-dec-det*: *is-dec-det* $d0$
 ⟨proof⟩

lemma *v0-code*[code]: $v0\text{-mpi}_b$ = *vec-to-bfun* ($\chi s. r\text{-min} / (1 - l)$)
 ⟨proof⟩

lemma *d0'-code*[code]: $d0'$ = *vec-to-fun* ($\chi s. (\text{LEAST } a. a \in A s)$)
 ⟨proof⟩

lemma *step-value-code*[code]: $L\text{-pow } v d m$ = ($L\text{-det } d \rightsquigarrow \text{Suc } m$) v
 ⟨proof⟩

partial-function (*tailrec*) *mpi-partial* **where** [*code*]: *mpi-partial eps*
 $d \ v \ m =$
 (let $d' = \text{policy-improvement } d \ v$ in (
 if $2 * l * \text{dist } v \ (\mathcal{L}_b \ v) < \text{eps} * (1 - l)$
 then (d', v)
 else *mpi-partial eps* $d' \ (L\text{-pow } v \ d' \ (m \ 0 \ v)) \ (\lambda n. m \ (\text{Suc } n))$))

lemma *mpi-partial-eq-algo*:
assumes $\text{eps} > 0 \ d \in D_D \ v \leq \mathcal{L}_b \ v$
shows *mpi-partial eps* $d \ v \ m = \text{mpi-algo } \text{eps} \ d \ v \ m$
 ⟨*proof*⟩

lift-definition *mpi-partial'* :: $\text{real} \Rightarrow ('s, 'a) \text{Fun} \Rightarrow ('s \Rightarrow_b \text{real}) \Rightarrow$
 $(\text{nat} \Rightarrow ('s \Rightarrow_b \text{real}) \Rightarrow \text{nat})$
 $\Rightarrow ('s, 'a) \text{Fun} \times ('s \Rightarrow_b \text{real})$ **is** *mpi-partial'*⟨*proof*⟩

lemma *mpi-partial'-code*[*code*]: *mpi-partial'* $\text{eps} \ d \ v \ m =$
 (let $d' = \text{policy-improvement}' \ d \ v$ in (
 if $2 * l * \text{dist } v \ (\mathcal{L}_b \ v) < \text{eps} * (1 - l)$
 then (d', v)
 else *mpi-partial'* $\text{eps} \ d' \ (L\text{-pow } v \ (\text{Rep-Fun } d') \ (m \ 0 \ v)) \ (\lambda n. m \ (\text{Suc } n))$))
 ⟨*proof*⟩

lemma *r-min-code*[*code-unfold*]: *r-min* = $(\text{MIN } s. \text{MIN } a. r(s, a))$
 ⟨*proof*⟩

lemma *mpi-user-code*[*code*]: *mpi-user eps* $m =$
 (if $\text{eps} \leq 0$ then *undefined* else
 let $(d, v) = \text{mpi-partial}' \ \text{eps} \ d0' \ v0\text{-mpi}_b \ m$ in $(\text{Rep-Fun } d, v)$)
 ⟨*proof*⟩
end

7.9 Auxiliary Equations

lemma [*code-unfold*]: $\text{dist} \ (f :: 'a :: \text{finite} \Rightarrow_b 'b :: \text{metric-space}) \ g = (\text{MAX}$
 $a. \text{dist} \ (\text{apply-bfun } f \ a) \ (g \ a))$
 ⟨*proof*⟩

lemma *member-code*[*code del*]: $x \in \text{List.coset } xs \iff \neg \text{List.member}$
 $xs \ x$
 ⟨*proof*⟩

lemma [*code*]: $\text{iarray-to-vec } v + \text{iarray-to-vec } u = (\text{Matrix-To-IArray.iarray-to-vec}$
 $(\text{Rep-iarray-type } v + \text{Rep-iarray-type } u))$
 ⟨*proof*⟩

lemma [*code*]: $\text{iarray-to-vec } v - \text{iarray-to-vec } u = (\text{Matrix-To-IArray.iarray-to-vec}$
 $(\text{Rep-iarray-type } v - \text{Rep-iarray-type } u))$

<proof>

lemma *matrix-to-iarray-minus*[code-unfold]: *matrix-to-iarray* ($A - B$)
= *matrix-to-iarray* $A - \text{matrix-to-iarray } B$
<proof>

declare *matrix-to-iarray-fst-Gauss-Jordan-PA*[code-unfold]

end
theory *Code-Mod*
imports *Code-DP*
begin

8 Code Generation for Concrete Finite MDPs

locale *mod-MDP* =
 fixes *transition* :: 's::{enum, mod-type} × 'a::{enum, mod-type} ⇒
 's pmf
 and *A* :: 's ⇒ 'a set
 and *reward* :: 's × 'a ⇒ real
 and *discount* :: real
begin

sublocale *mdp: vi-code*
 λs. (if Set.is-empty (A s) then UNIV else A s)
 transition
 reward
 (if $1 \leq \text{discount} \vee \text{discount} < 0$ then 0 else *discount*)
 defines $\mathcal{L}_b = \text{mdp}.\mathcal{L}_b$
 and *L-det* = *mdp.L-det*
 and *value-iteration* = *mdp.value-iteration*
 and *vi-policy'* = *mdp.vi-policy'*
 and *find-policy'* = *mdp.find-policy'*
 and *find-policy-impl* = *mdp.find-policy-impl*
 and *is-opt-act* = *mdp.is-opt-act*
 and *value-iteration-partial* = *mdp.value-iteration-partial*
 and *policy-iteration* = *mdp.policy-iteration*
 and *is-dec-det* = *mdp.is-dec-det*
 and *policy-step* = *mdp.policy-step*
 and *policy-improvement* = *mdp.policy-improvement*
 and *policy-eval* = *mdp.policy-eval*
 and *mk-markovian* = *mdp.mk-markovian*
 and *policy-eval'* = *mdp.policy-eval'*
 and *policy-iteration-partial'* = *mdp.policy-iteration-partial'*
 and *policy-iteration'* = *mdp.policy-iteration'*
 and *policy-iteration-policy-step'* = *mdp.policy-step'*
 and *policy-iteration-policy-eval'* = *mdp.policy-eval'*
and *policy-iteration-policy-improvement'* = *mdp.policy-improvement'*
 and *gs-iteration* = *mdp.gs-iteration*

```

and gs-iteration-partial = mdp.gs-iteration-partial
and vi-gs-policy = mdp.vi-gs-policy
and opt-policy-gs = mdp.opt-policy-gs
and opt-policy-gs'' = mdp.opt-policy-gs''
and P-mat = mdp.P-mat
and r-vec' = mdp.r-vec'
and GS-rec-funb = mdp.GS-rec-funb
and GS-iter-max = mdp.GS-iter-max
and GS-iter = mdp.GS-iter
and mpi-user = mdp.mpi-user
and v0-mpib = mdp.v0-mpib
and mpi-partial' = mdp.mpi-partial'
and L-pow = mdp.L-pow
and v0-mpi = mdp.v0-mpi
and r-min = mdp.r-min
and d0 = mdp.d0
and d0' = mdp.d0'
and νb = mdp.νb
and vi-test = mdp.vi-test
⟨proof⟩
end

```

```

global-interpretation mod-MDP transition A reward discount
for transition A reward discount
defines mod-MDP-ℒb = mdp.ℒb
and mod-MDP-ℒb-L-det = mdp.L-det
and mod-MDP-value-iteration = mdp.value-iteration
and mod-MDP-vi-policy' = mdp.vi-policy'
and mod-MDP-find-policy' = mdp.find-policy'
and mod-MDP-find-policy-impl = mdp.find-policy-impl
and mod-MDP-is-opt-act = mdp.is-opt-act
and mod-MDP-value-iteration-partial = mdp.value-iteration-partial
and mod-MDP-policy-iteration = mdp.policy-iteration
and mod-MDP-is-dec-det = mdp.is-dec-det
and mod-MDP-policy-step = mdp.policy-step
and mod-MDP-policy-improvement = mdp.policy-improvement
and mod-MDP-policy-eval = mdp.policy-eval
and mod-MDP-mk-markovian = mdp.mk-markovian
and mod-MDP-policy-eval' = mdp.policy-eval'
and mod-MDP-policy-iteration-partial' = mdp.policy-iteration-partial'
and mod-MDP-policy-iteration' = mdp.policy-iteration'
and mod-MDP-policy-iteration-policy-step' = mdp.policy-step'
and mod-MDP-policy-iteration-policy-eval' = mdp.policy-eval'
and mod-MDP-policy-iteration-policy-improvement' = mdp.policy-improvement'
and mod-MDP-gs-iteration = mdp.gs-iteration
and mod-MDP-gs-iteration-partial = mdp.gs-iteration-partial
and mod-MDP-vi-gs-policy = mdp.vi-gs-policy
and mod-MDP-opt-policy-gs = mdp.opt-policy-gs
and mod-MDP-opt-policy-gs'' = mdp.opt-policy-gs''

```

```

and mod-MDP-P-mat = mdp.P-mat
and mod-MDP-r-vec' = mdp.r-vec'
and mod-MDP-GS-rec-funb = mdp.GS-rec-funb
and mod-MDP-GS-iter-max = mdp.GS-iter-max
and mod-MDP-GS-iter = mdp.GS-iter
and mod-MDP-mpi-user = mdp.mpi-user
and mod-MDP-v0-mpib = mdp.v0-mpib
and mod-MDP-mpi-partial' = mdp.mpi-partial'
and mod-MDP-L-pow = mdp.L-pow
and mod-MDP-v0-mpi = mdp.v0-mpi
and mod-MDP-r-min = mdp.r-min
and mod-MDP-d0 = mdp.d0
and mod-MDP-d0' = mdp.d0'
and mod-MDP-νb = mdp.νb
and mod-MDP-vi-test = mdp.vi-test
⟨proof⟩

```

```

end
theory Code-Real-Approx-By-Float-Fix
imports
  HOL-Library.Code-Real-Approx-By-Float
  Gauss-Jordan.Code-Real-Approx-By-Float-Haskell
beginend

```

```

theory Code-Inventory
imports
  Code-Mod

  Code-Real-Approx-By-Float-Fix
begin

```

9 Inventory Management Example

```

lemma [code abstype]: embed-pmf (pmf P) = P
⟨proof⟩

```

```

lemmas [code-abbrev del] = pmf-integral-code-unfold

```

```

lemma [code-unfold]:
  measure-pmf.expectation P (f :: 'a :: enum ⇒ real) = (∑ x ∈ UNIV.
  pmf P x * f x)
⟨proof⟩

```

```

lemma [code]: pmf (return-pmf x) = (λy. indicat-real {y} x)
⟨proof⟩

```


lemma [code]:
 $pmf (bind\text{-}pmf\ N\ f) = (\lambda i :: 'a. measure\text{-}pmf.expectation\ N\ (\lambda(x :: 'b :: enum). pmf\ (f\ x)\ i))$
 ⟨proof⟩

lemma *pmf-finite-le*: $finite\ (X :: ('a::finite)\ set) \implies sum\ (pmf\ p)\ X \leq 1$
 ⟨proof⟩

lemma *mod-less-diff*:
assumes $0 < (x :: 's :: \{mod\text{-}type\})\ x \leq y$
shows $y - x < y$
 ⟨proof⟩

locale *inventory* =
fixes *fixed-cost* :: real
and *var-cost* :: '*s*::{*mod-type*, *finite*} \Rightarrow real
and *inv-cost* :: '*s* \Rightarrow real
and *demand-prob* :: '*s* pmf
and *revenue* :: '*s* \Rightarrow real
and *discount* :: real
begin
definition *order-cost* *u* = (if *u* = 0 then 0 else *fixed-cost* + *var-cost* *u*)
definition *prob-ge-inv* *u* = $1 - (\sum j < u. pmf\ demand\text{-}prob\ j)$
definition *exp-rev* *u* = $(\sum j < u. revenue\ j * pmf\ demand\text{-}prob\ j) + revenue\ u * prob\text{-}ge\text{-}inv\ u$
definition *reward* *sa* = (case *sa* of (*s*,*a*) \Rightarrow *exp-rev* (*s* + *a*) - *order-cost* *a* - *inv-cost* (*s* + *a*))
lift-definition *transition* :: ('*s* × '*s*) \Rightarrow '*s* pmf **is** $\lambda(s, a)\ s'$.
 (if $CARD('s) \leq Rep\ s + Rep\ a$
 then (if $s' = 0$ then 1 else 0)
 else (if $s + a < s'$ then 0 else
 if $s' = 0$ then *prob-ge-inv* (*s*+*a*)
 else *pmf* *demand-prob* (*s* + *a* - *s'*)))

⟨proof⟩

definition *A-inv* (*s*::'*s*) = {*a*::'*s*. $Rep\ s + Rep\ a < CARD('s)$ }

end

definition *var-cost-lin* (*c*::real) *n* = $c * Rep\ n$

definition *inv-cost-lin* (*c*::real) *n* = $c * Rep\ n$

definition *revenue-lin* (*c*::real) *n* = $c * Rep\ n$

lift-definition *demand-unif* :: ('*a*::*finite*) pmf **is** $\lambda. 1 / card\ (UNIV::'a\ set)$
 ⟨proof⟩

lift-definition *demand-three* :: \mathcal{P} pmf is $\lambda i.$ if $i = 1$ then $1/4$ else if $i = 2$ then $1/2$ else $1/4$
 ⟨proof⟩

abbreviation *fixed-cost* $\equiv 4$
abbreviation *var-cost* \equiv *var-cost-lin 2*
abbreviation *inv-cost* \equiv *inv-cost-lin 1*
abbreviation *revenue* \equiv *revenue-lin 8*
abbreviation *discount* $\equiv 0.99$
type-synonym *capacity* = 30

lemma *card-ge-2-imp-ne*: $CARD('a) \geq 2 \implies \exists (x::'a::finite) y::'a. x \neq y$
 ⟨proof⟩

global-interpretation *inventory-ex*: *inventory fixed-cost var-cost:: capacity* \Rightarrow *real inv-cost demand-unif revenue discount*
defines *A-inv* = *inventory-ex.A-inv*
and *transition* = *inventory-ex.transition*
and *reward* = *inventory-ex.reward*
and *prob-ge-inv* = *inventory-ex.prob-ge-inv*
and *order-cost* = *inventory-ex.order-cost*
and *exp-rev* = *inventory-ex.exp-rev*⟨proof⟩

abbreviation *K* \equiv *inventory-ex.transition*
abbreviation *A* \equiv *inventory-ex.A-inv*
abbreviation *r* \equiv *inventory-ex.reward*
abbreviation *l* $\equiv 0.95$
definition *eps* = 0.1

definition *fun-to-list* $f = \text{map } f \text{ (sorted-list-of-set UNIV)}$
definition *benchmark-gs* ($- :: \text{unit}$) = $\text{map Rep (fun-to-list (vi-policy' K A r l eps 0))}$
definition *benchmark-vi* ($- :: \text{unit}$) = $\text{map Rep (fun-to-list (vi-gs-policy K A r l eps 0))}$
definition *benchmark-mpi* ($- :: \text{unit}$) = $\text{map Rep (fun-to-list (fst (mpi-user K A r l eps (\lambda -. \mathcal{P}))))}$
definition *benchmark-pi* ($- :: \text{unit}$) = $\text{map Rep (fun-to-list (policy-iteration K A r l 0))}$

fun *vs-n* **where**
vs-n 0 $v = v$
 | *vs-n (Suc n)* $v = \text{vs-n } n \text{ (mod-MDP-}\mathcal{L}_b \text{ K A r l } v)$

definition *vs-n'* $n = \text{vs-n } n \text{ 0}$

definition *benchmark-vi-n* $n = \text{(fun-to-list (vs-n } n \text{ 0))}$
definition *benchmark-vi-nopol* = $\text{(fun-to-list (mod-MDP-value-iteration$

K A r l (1/10) 0)

export-code *dist vs-n' benchmark-vi-nopol benchmark-vi-n nat-of-integer
integer-of-int benchmark-gs benchmark-vi benchmark-mpi benchmark-pi*
in Haskell module-name DP

export-code *integer-of-int benchmark-gs benchmark-vi benchmark-mpi
benchmark-pi* **in SML module-name DP**

end

theory *Code-Gridworld*
imports
 Code-Mod
begin

10 Gridworld Example

lemma [*code abstype*]: *embed-pmf (pmf P) = P*
 <proof>

lemmas [*code-abbrev del*] = *pmf-integral-code-unfold*

lemma [*code-unfold*]:
 measure-pmf.expectation P (f :: 'a :: enum \Rightarrow real) = ($\sum x \in UNIV.$
*pmf P x * f x)*
 <proof>

lemma [*code*]: *pmf (return-pmf x) = ($\lambda y.$ *indicat-real {y} x*)*
 <proof>

lemma [*code*]:
 *pmf (bind-pmf N f) = ($\lambda i :: 'a.$ *measure-pmf.expectation N ($\lambda(x ::$*
*'b :: enum). *pmf (f x) i)*)*
 *<proof>**

type-synonym *state-robot = 13*

definition *from-state x = (Rep x div 4, Rep x mod 4)*

definition *to-state x = (Abs (fst x * 4 + snd x) :: state-robot)*

type-synonym *action-robot = 4*

fun *A-robot* :: *state-robot* \Rightarrow *action-robot set* **where**
A-robot pos = *UNIV*

abbreviation *noise* \equiv (*0.2* :: *real*)

lift-definition *add-noise* :: *action-robot* \Rightarrow *action-robot pmf* **is** λdet
rnd. (
 if det = rnd then 1 - noise else if det = rnd - 1 \vee det = rnd + 1
 then noise / 2 else 0)
\langle proof \rangle

fun *r-robot* :: (*state-robot* \times *action-robot*) \Rightarrow *real* **where**
r-robot (s,a) = (
 if from-state s = (2,3) then 1 else
 if from-state s = (1,3) then -1 else
 if from-state s = (3,0) then 0 else
 0)

fun *K-robot* :: (*state-robot* \times *action-robot*) \Rightarrow *state-robot pmf* **where**
K-robot (loc, a) =
 do {
 a \leftarrow *add-noise a*;
 let (y, x) = from-state loc;
 let (y', x') =
 (*if a = 0 then (y + 1, x)*
 else if a = 1 then (y, x+1)
 else if a = 2 then (y-1, x)
 else if a = 3 then (y, x-1)
 else undefined);
 return-pmf (
 if (y,x) = (2,3) \vee (y,x) = (1,3) \vee (y,x) = (3,0)
 then to-state (3,0)
 else if y' < 0 \vee y' > 2 \vee x' < 0 \vee x' > 3 \vee (y',x') = (1,1)
 then to-state (y, x)
 else to-state (y', x'))
 }

definition *l-robot* = *0.9*

lemma *vi-code A-robot r-robot l-robot*
\langle proof \rangle

abbreviation *to-gridworld f* \equiv *f K-robot r-robot l-robot*

abbreviation *to-gridworld' f* \equiv *f K-robot A-robot r-robot l-robot*

abbreviation *gridworld-policy-eval'* \equiv *to-gridworld mod-MDP-policy-eval'*

abbreviation *gridworld-policy-step'* \equiv *to-gridworld' mod-MDP-policy-iteration-policy-step'*

abbreviation *gridworld-mpi-user* \equiv *to-gridworld' mod-MDP-mpi-user*

abbreviation *gridworld-opt-policy-gs* \equiv *to-gridworld' mod-MDP-opt-policy-gs*
abbreviation *gridworld- \mathcal{L}_b* \equiv *to-gridworld' mod-MDP- \mathcal{L}_b*
abbreviation *gridworld-find-policy'* \equiv *to-gridworld' mod-MDP-find-policy'*
abbreviation *gridworld-GS-rec-fun $_b$* \equiv *to-gridworld' mod-MDP-GS-rec-fun $_b$*
abbreviation *gridworld-vi-policy'* \equiv *to-gridworld' mod-MDP-vi-policy'*
abbreviation *gridworld-vi-gs-policy* \equiv *to-gridworld' mod-MDP-vi-gs-policy*
abbreviation *gridworld-policy-iteration* \equiv *to-gridworld' mod-MDP-policy-iteration*

fun *pi-robot-n* **where**

pi-robot-n 0 *d* = (*d*, *gridworld-policy-eval' d*) |
pi-robot-n (*Suc n*) *d* = *pi-robot-n n* (*gridworld-policy-step' d*)

definition *mpi-robot eps* = *gridworld-mpi-user eps* (λ -. 3)

fun *gs-robot-n* **where**

gs-robot-n (0 :: *nat*) *v* = (*gridworld-opt-policy-gs v*, *v*) |
gs-robot-n (*Suc n* :: *nat*) *v* = *gs-robot-n n* (*gridworld-GS-rec-fun $_b$ v*)

fun *vi-robot-n* **where**

vi-robot-n (0 :: *nat*) *v* = (*gridworld-find-policy' v*, *v*) |
vi-robot-n (*Suc n* :: *nat*) *v* = *vi-robot-n n* (*gridworld- \mathcal{L}_b v*)

definition *mpi-result eps* =

(*let* (*d*, *v*) = *mpi-robot eps* *in* (*d,v*))

definition *gs-result n* =

(*let* (*d,v*) = *gs-robot-n n* 0 *in* (*d,v*))

definition *vi-result-n n* =

(*let* (*d*, *v*) = *vi-robot-n n* 0 *in* (*d,v*))

definition *pi-result-n n* =

(*let* (*d*, *v*) = *pi-robot-n n* (*vec-to-fun* 0) *in* (*d,v*))

definition *fun-to-list f* = *map f* (*sorted-list-of-set UNIV*)

definition *benchmark-gs* = *fun-to-list* (*gridworld-vi-policy' 0.1* 0)

definition *benchmark-vi* = *fun-to-list* (*gridworld-vi-gs-policy 0.1* 0)

definition *benchmark-mpi* = *fun-to-list* (*fst* (*gridworld-mpi-user 0.1* (λ - -. 3)))

definition *benchmark-pi* = *fun-to-list* (*gridworld-policy-iteration* 0)

export-code *benchmark-gs benchmark-vi benchmark-mpi benchmark-pi*

in *Haskell* **module-name** *DP*

export-code *benchmark-gs benchmark-vi benchmark-mpi benchmark-pi*

in *SML* **module-name** *DP*

end

```
theory Examples  
imports  
  Code-Inventory  
  Code-Gridworld  
begin  
end
```

References

- [1] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley, 1994.