

Markov Decision Processes with Rewards

Maximilian Schäffeler and Mohammad Abdulaziz

December 28, 2021

Abstract

We present a formalization of Markov Decision Processes with rewards. In particular we first build on Hölzl’s formalization [1] of MDPs and extend them with rewards. We proceed with an analysis of the expected total discounted reward criterion for infinite horizon MDPs. The central result is the construction of the iteration rule for the Bellman operator. We prove the optimality equations for this operator and show the existence of an optimal stationary deterministic solution. The analysis can be used to obtain dynamic programming algorithms such as value iteration and policy iteration to solve Markov Decision Processes with formal guarantees. Our formalization is based upon chapters 5 and 6 in Puterman’s book [2].

Contents

1	Bounded Functions	3
1.1	Definition	3
1.2	Supremum Norm	5
1.3	Complete Space	7
1.4	Order Instance	8
1.5	Miscellaneous	9
1.6	Bounded Functions and Vectors	9
2	Bounded Linear Functions	10
2.1	Composition	10
2.2	Power	10
2.3	Geometric Sum	11
2.4	Inverses	12
2.5	Norm	14
2.6	Miscellaneous	15
3	Auxiliary Lemmas	16
3.1	Summability	16
3.2	Infinite sums	16
3.3	Bounded Functions	16
3.4	Push-Forward of a Bounded Function	16
3.5	Boundedness	17

3.6	Probability Theory	17
3.7	Argmax	18
3.8	Contraction Mappings	19
3.9	Limits	20
3.10	Supremum	20
4	Discrete-Time Markov Decision Processes with Arbitrary State Spaces	21
4.1	Definition and Basic Properties	21
4.2	Policies	23
4.3	Successor Policy	25
4.4	Single-Step Distribution	25
4.5	Initial State-Action Distribution	26
4.6	Sequence Space of the MDP	27
4.7	Measurability of the Sequence Space	27
4.8	Iteration Rule	28
4.9	Stream Space of the MDP	28
5	Markov Decision Processes with Discrete State Spaces	30
5.1	Policies	31
5.1.1	Successor Policy	33
5.2	Stream Space of the MDP	33
5.2.1	Initial State-Action Distribution	33
5.2.2	Decomposition of the Stream Space	35
5.2.3	A Denotational View on the Stochastic Process	36
5.2.4	State Process	37
5.2.5	The Conditional Distribution of Actions	37
5.2.6	Action Process	38
5.3	Restriction to Markovian Policies	39
5.4	MDPs without Initial Distribution	39
6	Markov Decision Processes with Rewards	42
6.1	Basic Properties	43
6.2	Summability	43
6.3	Reward over a Trace	44
6.4	Integrals over Rewards	44
6.5	Expected Total Discounted Reward	44
6.6	Expected Finite-Horizon Discounted Reward	44
6.6.1	Optimal Reward	45
6.7	Reward of a Decision Rule	46
6.8	Push-Forward of a Function Through the MDP	46
6.9	The Bellman Operator L	51
6.10	Optimality Equations	52
6.11	Monotonicity	52
6.12	Properties of Solutions of the Optimality Equations	54
6.13	Solutions to the Optimality Equation	54
6.13.1	\mathcal{L}_b and L are Contraction Mappings	54
6.13.2	Existence of a Fixpoint of \mathcal{L}_b	54
6.14	Existence of Optimal Policies	55
6.14.1	Conserving Decision Rules are Optimal	56

6.14.2	Bellman Operator for Single Actions	57
6.14.3	Equivalences involving \mathcal{L}_b	57
6.14.4	Deterministic Decision Rules are Optimal	58
6.14.5	Optimal Decision Rules for Finite Action Spaces	58
6.14.6	Existence of Epsilon-Optimal Policies	58
6.15	More Restrictive MDP Locales	59

1 Bounded Functions

theory *Bounded-Functions*

imports

HOL.Topological-Spaces

HOL-Analysis.Uniform-Limit

HOL-Probability.Probability

begin

1.1 Definition

definition *bfun* = {*f*. *bounded* (*range* *f*)}

typedef (**overloaded**) (*'a*, *'b*) *bfun* ((*-* \Rightarrow_b *-*) [22] 21) =

bfun::(*'a* \Rightarrow *'b* :: *metric-space*) *set*

morphisms *apply-bfun* *Bfun*

<proof>

declare [[*coercion* *apply-bfun* :: (*'a* \Rightarrow_b (*'b* :: *metric-space*)) \Rightarrow *'a* \Rightarrow *'b*]]

setup-lifting *type-definition-bfun*

lemma *bounded-apply-bfun*[*intro*, *simp*]: *bounded* (*range* (*apply-bfun* *x*))

<proof>

lemma *bfun-eqI*[*intro*]: ($\bigwedge x.$ *apply-bfun* *f* *x* = *apply-bfun* *g* *x*) \Longrightarrow *f* = *g*

<proof>

lemma *bfun-eqD*[*dest*]: *f* = *g* \Longrightarrow ($\bigwedge x.$ *apply-bfun* *f* *x* = *apply-bfun* *g* *x*)

<proof>

lemma *bfunE*:

assumes *f* \in *bfun*

obtains *g* **where** *f* = *apply-bfun* *g*

<proof>

lemma *const-bfun*: ($\lambda x.$ *b*) \in *bfun*

<proof>

lift-definition *const-bfun*::'b \Rightarrow ('a \Rightarrow_b ('b :: *metric-space*)) **is** $\lambda(c::'b)$
 \cdot . c
 <proof>

lemma *bounded-dist-le-SUP-dist*:
 bounded (range f) \Longrightarrow bounded (range g) \Longrightarrow dist (f x) (g x) \leq (SUP
 x. dist (f x) (g x))
 <proof>

instantiation *bfun* :: (*type*, *metric-space*) *metric-space*
begin

lift-definition *dist-bfun* :: ('a \Rightarrow_b 'b) \Rightarrow ('a \Rightarrow_b 'b) \Rightarrow *real*
is $\lambda f g$. (SUP x. dist (f x) (g x)) <proof>

definition *uniformity-bfun* :: (('a \Rightarrow_b 'b) \times 'a \Rightarrow_b 'b) *filter*
where *uniformity-bfun* = (INF e \in {0 <.. ∞ }. principal {(x, y). dist x
 y < e})

definition *open-bfun* :: ('a \Rightarrow_b 'b) *set* \Rightarrow *bool*
where *open-bfun* S = ($\forall x \in S$. $\forall_F (x', y)$ in *uniformity*. $x' = x \longrightarrow$
 y \in S)

lemma *dist-bounded*:
fixes f g :: 'a \Rightarrow_b 'b
shows dist (f x) (g x) \leq dist f g
 <proof>

lemma *dist-bound*:
fixes f g :: 'a \Rightarrow_b ('b :: *metric-space*)
assumes $\bigwedge x$. dist (f x) (g x) \leq b
shows dist f g \leq b
 <proof>

lemma *dist-fun-lt-imp-dist-val-lt*:
fixes f g :: 'a \Rightarrow_b 'b
assumes dist f g < e
shows dist (f x) (g x) < e
 <proof>

instance
 <proof>

end

lift-definition *PiC*::'a *set* \Rightarrow ('a \Rightarrow ('b :: *metric-space*) *set*) \Rightarrow ('a
 \Rightarrow_b 'b) *set*

is $\lambda I X. Pi I X \cap bfun$
<proof>

lemma *mem-PiC-iff*: $x \in PiC I X \longleftrightarrow apply-bfun x \in Pi I X$
<proof>

lemmas *mem-PiCD = mem-PiC-iff*[*THEN iffD1*]
and *mem-PiCI = mem-PiC-iff*[*THEN iffD2*]

lemma *tendsto-bfun-uniform-limit*:
fixes $f::'i \Rightarrow 'a \Rightarrow_b ('b :: metric-space)$
assumes $(f \longrightarrow l) F$
shows *uniform-limit UNIV f l F*
<proof>

lemma *uniform-limit-tendsto-bfun*:
fixes $f::'i \Rightarrow 'a \Rightarrow_b ('b :: metric-space)$
and $l::'a \Rightarrow_b 'b$
assumes *uniform-limit UNIV f l F*
shows $(f \longrightarrow l) F$
<proof>

1.2 Supremum Norm

instantiation *bfun* :: $(type, real-normed-vector) real-vector$
begin

lemma *uminus-cont*: $f \in bfun \implies (\lambda x. - f x) \in bfun$ **for** $f::'a \Rightarrow 'b$
<proof>

lemma *plus-cont*: $f \in bfun \implies g \in bfun \implies (\lambda x. f x + g x) \in bfun$
for $f g::'a \Rightarrow 'b$
<proof>

lemma *minus-cont*: $f \in bfun \implies g \in bfun \implies (\lambda x. f x - g x) \in bfun$
for $f g::'a \Rightarrow 'b$
<proof>

lemma *scaleR-cont*: $f \in bfun \implies (\lambda x. a *_R f x) \in bfun$ **for** $f::'a \Rightarrow 'b$
<proof>

lemma *bfun-normI*[*intro*]: $(\bigwedge x. norm (f x) \leq b) \implies f \in bfun$
<proof>

lift-definition *uminus-bfun*:: $('a \Rightarrow_b 'b) \Rightarrow ('a \Rightarrow_b 'b)$ **is** $\lambda f x. - f x$
<proof>

lift-definition *plus-bfun*:: $('a \Rightarrow_b 'b) \Rightarrow ('a \Rightarrow_b 'b) \Rightarrow 'a \Rightarrow_b 'b$ **is** λf

$g\ x.\ f\ x + g\ x$
 $\langle proof \rangle$

lift-definition $minus-bfun::('a \Rightarrow_b 'b) \Rightarrow ('a \Rightarrow_b 'b) \Rightarrow 'a \Rightarrow_b 'b$ is
 $\lambda f\ g\ x.\ f\ x - g\ x$
 $\langle proof \rangle$

lift-definition $zero-bfun::'a \Rightarrow_b 'b$ is $\lambda-. 0$
 $\langle proof \rangle$

lemma $const-bfun-0-eq-0[simp]: const-bfun\ 0 = 0$
 $\langle proof \rangle$

lift-definition $scaleR-bfun::real \Rightarrow ('a \Rightarrow_b 'b) \Rightarrow 'a \Rightarrow_b 'b$ is $\lambda r\ g\ x.$
 $r *_{\mathbb{R}} g\ x$
 $\langle proof \rangle$

lemmas $[simp] =$
 $const-bfun.rep-eq$
 $uminus-bfun.rep-eq$
 $plus-bfun.rep-eq$
 $minus-bfun.rep-eq$
 $zero-bfun.rep-eq$
 $scaleR-bfun.rep-eq$

instance
 $\langle proof \rangle$
end

lemma $scaleR-cont': f \in bfun \Longrightarrow (\lambda x.\ a * f\ x) \in bfun$ for $f :: 'a \Rightarrow$
 $real$
 $\langle proof \rangle$

lemma $bfun-norm-le-SUP-norm:$
 $f \in bfun \Longrightarrow norm\ (f\ x) \leq (SUP\ x.\ norm\ (f\ x))$
 $\langle proof \rangle$

instantiation $bfun :: (type, real-normed-vector) real-normed-vector$
begin

definition $norm-bfun :: ('a, 'b) bfun \Rightarrow real$
where $norm-bfun\ f = dist\ f\ 0$

definition $sgn\ (f::('a,'b) bfun) = f /_{\mathbb{R}} norm\ f$

instance
 $\langle proof \rangle$
end

lemma *norm-bfun-def'*: $\text{norm } f = (\bigsqcup x. \text{norm } ((f :: 'a \Rightarrow_b 'b :: \text{real-normed-vector}) x))$
 ⟨proof⟩

lemma *norm-le-norm-bfun*: $\text{norm } (\text{apply-bfun } f x) \leq \text{norm } f$
 ⟨proof⟩

lemma *abs-le-norm-bfun*: $\text{abs } (\text{apply-bfun } f x) \leq \text{norm } f$
 ⟨proof⟩

lemma *le-norm-bfun*: $\text{apply-bfun } f x \leq \text{norm } f$
 ⟨proof⟩

1.3 Complete Space

lemma *tendsto-add*: $P \longrightarrow (L :: 'a :: \text{real-normed-vector}) \Longrightarrow (\lambda n. P n + c) \longrightarrow L + c$
 ⟨proof⟩

lemma *lim-add*: $\text{convergent } P \Longrightarrow \text{lim } (\lambda n. P n + (c :: 'a :: \text{real-normed-vector})) = \text{lim } P + c$
 ⟨proof⟩

lemma *complete-bfun*:
assumes *cauchy-f*: $\text{Cauchy } (f :: \text{nat} \Rightarrow ('a, 'b :: \{\text{complete-space, real-normed-vector}\}) \text{ bfun})$
shows *convergent f*
 ⟨proof⟩

lemma *norm-bound*:
fixes $f :: ('a, 'b :: \text{real-normed-vector}) \text{ bfun}$
assumes $\bigwedge x. \text{norm } (\text{apply-bfun } f x) \leq b$
shows $\text{norm } f \leq b$
 ⟨proof⟩

lemma *bfun-bounded-norm-range*: $\text{bounded } (\text{range } (\lambda s. \text{norm } (\text{apply-bfun } v s)))$
 ⟨proof⟩

instance *bfun :: (type, banach) banach*
 ⟨proof⟩

lemma *bfun-prob-space-integrable*:
assumes *prob-space S v* $\in \text{borel-measurable } S$
assumes $(v :: 'a \Rightarrow 'b :: \{\text{second-countable-topology, banach}\}) \in \text{bfun}$
shows *integrable S v*
 ⟨proof⟩

lemma *bfun-integral-bound*:
assumes $(v :: 'a \Rightarrow 'c :: \{\text{euclidean-space}\}) \in \text{bfun}$
shows $(\lambda S. \int x. v x \partial(S :: 'a \text{ pmf})) \in \text{bfun}$
 $\langle \text{proof} \rangle$

lemma *scale-bfun[intro]*: $f \in \text{bfun} \implies (\lambda x. (k :: \text{real}) * f x) \in \text{bfun}$
 $\langle \text{proof} \rangle$

lemma *bfun-spec[intro]*: $f \in \text{bfun} \implies (\lambda x. f (g x)) \in \text{bfun}$
 $\langle \text{proof} \rangle$

lemma *apply-bfun-bfun[simp]*: $\text{apply-bfun } f \in \text{bfun}$
 $\langle \text{proof} \rangle$

lemma *bfun-integral-bound'[intro]*: $(v :: 'a \Rightarrow 'c :: \{\text{euclidean-space}\}) \in \text{bfun} \implies$
 $(\lambda S. \int x. v x \partial((F S) :: 'a \text{ pmf})) \in \text{bfun}$
 $\langle \text{proof} \rangle$

lift-definition *bfun-comp* :: $('a \Rightarrow 'b) \Rightarrow ('b \Rightarrow_b 'c :: \text{metric-space}) \Rightarrow$
 $('a \Rightarrow_b 'c)$ **is**
 $\lambda g \text{ bf } x. \text{bf } (g x)$
 $\langle \text{proof} \rangle$

1.4 Order Instance

class *ordered-real-normed-vector* = *real-normed-vector* + *ordered-real-vector*

instance *real* :: *ordered-real-normed-vector*
 $\langle \text{proof} \rangle$

instantiation *bfun* :: $(-, \text{ordered-real-normed-vector}) \text{ ordered-real-normed-vector}$
begin

definition *less-eq-bfun* $f g \equiv \forall x. \text{apply-bfun } f x \leq \text{apply-bfun } g x$

definition *less-bfun* $f g \equiv \forall x. \text{apply-bfun } f x \leq \text{apply-bfun } g x \wedge (\exists y. f y < g y)$

instance
 $\langle \text{proof} \rangle$
end

lemma *less-eq-bfunI[intro]*: $(\wedge x. \text{apply-bfun } f x \leq \text{apply-bfun } g x) \implies f \leq g$
 $\langle \text{proof} \rangle$

lemma *less-eq-bfunD[dest]*: $f \leq g \implies (\wedge x. \text{apply-bfun } f x \leq \text{apply-bfun } g x)$

⟨proof⟩

1.5 Miscellaneous

instantiation *bfun* :: (*type*, *one*) *one* **begin**

lift-definition *one-bfun* :: '*s* ⇒_{*b*} *real* **is** λ*x*. 1
⟨proof⟩

instance
⟨proof⟩
end

declare *one-bfun.rep-eq* [*simp*]

lemma *apply-bfun-one* [*simp*]: *apply-bfun* (1 :: - ⇒_{*b*} *real*) *x* = 1
⟨proof⟩

lemma *norm-bfun-one*[*simp*]: *norm* (1 :: '*a* ⇒_{*b*} *real*) = 1
⟨proof⟩

lemma *range-bfunI*[*intro*]: *bounded* (*range* *f*) ⇒ *f* ∈ *bfun*
⟨proof⟩

lemma *finite-bfun*[*simp*]: (λ(*i*::::finite). *f* *i*) ∈ *bfun*
⟨proof⟩

lemma *bounded-apply-bfun'*:
 assumes *bounded* ((*F* :: '*c* ⇒ '*d* ⇒_{*b*} '*b*::*real-normed-vector*) '*S*)
 shows *bounded* ((λ*b*. (*F* *b*) *x*) '*S*)
⟨proof⟩

lemma *bfun-tendsto-apply-bfun*:
 assumes *h*: (*F* :: (*nat* ⇒ '*a* ⇒_{*b*} *real*)) ⟶ (y :: '*a* ⇒_{*b*} *real*)
 shows (λ*n*. *F* *n* *x*) ⟶ *y* *x*
⟨proof⟩

1.6 Bounded Functions and Vectors

lemma *vec-bfun*[*simp*, *intro*]: (\$) *x* ∈ *bfun*
⟨proof⟩

lemma *norm-bfun-le-norm-vec*: *norm* (*bfun.Bfun* ((*\$*) (*x* :: *real*^{*c*} ::
finite))) ≤ *norm* *x*
⟨proof⟩

lemma *bounded-linear-bfun-nth*: *bounded-linear* *f* ⇒ *bounded-linear*
(λ*v*. *bfun.Bfun* ((*\$*) (*f* *v*)))

<proof>

lemma *norm-vec-le-norm-bfun*:

$\text{norm } (\text{vec-lambda } (\text{apply-bfun } (x :: 'd::\text{finite} \Rightarrow_b \text{real}))) \leq \text{norm } x * \text{card } (\text{UNIV} :: 'd \text{ set})$

<proof>

end

2 Bounded Linear Functions

theory *Blinfun-Util*

imports

HOL-Analysis.Bounded-Linear-Function

Bounded-Functions

begin

2.1 Composition

lemma *blinfun-compose-id[simp]*:

$\text{id-blinfun } o_L f = f$

$f o_L \text{id-blinfun} = f$

<proof>

lemma *blinfun-compose-assoc*: $F o_L G o_L H = F o_L (G o_L H)$

<proof>

lemma *blinfun-compose-diff-right*: $f o_L (g - h) = (f o_L g) - (f o_L h)$

<proof>

2.2 Power

overloading

$\text{blinfunpow} \equiv \text{compow} :: \text{nat} \Rightarrow ('a::\text{real-normed-vector} \Rightarrow_L 'a) \Rightarrow ('a \Rightarrow_L 'a)$

begin

primrec *blinfunpow* :: $\text{nat} \Rightarrow ('a::\text{real-normed-vector} \Rightarrow_L 'a) \Rightarrow ('a \Rightarrow_L 'a)$

where

$\text{blinfunpow } 0 f = \text{id-blinfun}$

| $\text{blinfunpow } (\text{Suc } n) f = f o_L \text{blinfunpow } n f$

end

lemma *bounded-pow-blinfun[intro]*:

assumes *bounded* ($\text{range } (F::\text{nat} \Rightarrow 'a::\text{real-normed-vector} \Rightarrow_L 'a)$)

shows *bounded* ($\text{range } (\lambda t. (F t) \wedge\wedge (\text{Suc } n))$)

<proof>

lemma *blincomp-scaleR-right*: $(a *_R (F :: 'a :: \text{real-normed-vector} \Rightarrow_L 'a)) \widehat{\sim} t = a \widehat{\sim} t *_R F \widehat{\sim} t$
 ⟨proof⟩

lemma *summable-inv-Q*:
fixes $Q :: 'a :: \text{banach} \Rightarrow_L 'a$
assumes *onorm-le*: $\text{norm } (id\text{-blinfun} - Q) < 1$
shows *summable* $(\lambda n. (id\text{-blinfun} - Q) \widehat{\sim} n)$
 ⟨proof⟩

lemma *blinfunpow-assoc*: $(F :: 'a :: \text{real-normed-vector} \Rightarrow_L 'a) \widehat{\sim} (Suc\ n) = (F \widehat{\sim} n) \circ_L F$
 ⟨proof⟩

lemma *norm-blinfunpow-le*: $\text{norm } ((f :: 'b :: \text{real-normed-vector} \Rightarrow_L 'b) \widehat{\sim} n) \leq \text{norm } f \widehat{\sim} n$
 ⟨proof⟩

lemma *blinfunpow-nonneg*:
assumes $\bigwedge v. 0 \leq v \implies 0 \leq \text{blinfun-apply } (f :: ('b :: \{\text{ord, real-normed-vector}\} \Rightarrow_L 'b)) v$
shows $0 \leq v \implies 0 \leq (f \widehat{\sim} n) v$
 ⟨proof⟩

lemma *blinfunpow-mono*:
assumes $\bigwedge u v. u \leq v \implies (f :: 'b :: \{\text{ord, real-normed-vector}\} \Rightarrow_L 'b) u \leq f v$
shows $u \leq v \implies (f \widehat{\sim} n) u \leq (f \widehat{\sim} n) v$
 ⟨proof⟩

lemma *banach-blinfun*:
fixes $C :: 'b :: \{\text{real-normed-vector, complete-space}\} \Rightarrow_L 'b$
assumes *norm C < 1*
shows $\exists! v. C v = v \bigwedge v. (\lambda n. (C \widehat{\sim} n) v) \longrightarrow (THE v. C v = v)$
 ⟨proof⟩

2.3 Geometric Sum

lemma *inv-one-sub-Q*:
fixes $Q :: 'a :: \text{banach} \Rightarrow_L 'a$
assumes *onorm-le*: $\text{norm } (id\text{-blinfun} - Q) < 1$
shows $(Q \circ_L (\sum i. (id\text{-blinfun} - Q) \widehat{\sim} i)) = id\text{-blinfun}$
and $(\sum i. (id\text{-blinfun} - Q) \widehat{\sim} i) \circ_L Q = id\text{-blinfun}$
 ⟨proof⟩

lemma *inv-norm-le*:
fixes $Q :: 'a :: \text{banach} \Rightarrow_L 'a$
assumes *norm Q < 1*
shows $(id\text{-blinfun} - Q) \circ_L (\sum i. Q \widehat{\sim} i) = id\text{-blinfun}$

$(\sum i. Q \tilde{i}) \circ_L (id\text{-blinfun} - Q) = id\text{-blinfun}$
 ⟨proof⟩

lemma *inv-norm-le'*:
fixes $Q :: 'a :: \text{banach} \Rightarrow_L 'a$
assumes $\text{norm } Q < 1$
shows $(id\text{-blinfun} - Q) ((\sum i. Q \tilde{i}) x) = x$
 $(\sum i. Q \tilde{i}) ((id\text{-blinfun} - Q) x) = x$
 ⟨proof⟩

2.4 Inverses

definition $is\text{-inverse}_L X Y \longleftrightarrow X \circ_L Y = id\text{-blinfun} \wedge Y \circ_L X = id\text{-blinfun}$

abbreviation $invertible_L X \equiv \exists X'. is\text{-inverse}_L X X'$

lemma *is-inverse_L-I[intro]*:
assumes $X \circ_L Y = id\text{-blinfun}$ $Y \circ_L X = id\text{-blinfun}$
shows $is\text{-inverse}_L X Y$
 ⟨proof⟩

lemma *is-inverse_L-D[dest]*:
assumes $is\text{-inverse}_L X Y$
shows $X \circ_L Y = id\text{-blinfun}$ $Y \circ_L X = id\text{-blinfun}$
 ⟨proof⟩

lemma *invertible_L-D[dest]*:
assumes $invertible_L f$
obtains g **where** $f \circ_L g = id\text{-blinfun}$ $g \circ_L f = id\text{-blinfun}$
 ⟨proof⟩

lemma *invertible_L-I[intro]*:
assumes $f \circ_L g = id\text{-blinfun}$ $g \circ_L f = id\text{-blinfun}$
shows $invertible_L f$
 ⟨proof⟩

lemma *is-inverse_L-comm*: $is\text{-inverse}_L X Y \longleftrightarrow is\text{-inverse}_L Y X$
 ⟨proof⟩

lemma *is-inverse_L-unique*: $is\text{-inverse}_L f g \Longrightarrow is\text{-inverse}_L f h \Longrightarrow g = h$
 ⟨proof⟩

lemma *is-inverse_L-ex1*: $is\text{-inverse}_L f g \Longrightarrow \exists! h. is\text{-inverse}_L f h$
 ⟨proof⟩

lemma *is-inverse_L-ex1'*: $\exists x. is\text{-inverse}_L f x \Longrightarrow \exists! x. is\text{-inverse}_L f x$
 ⟨proof⟩

definition $inv_L f = (THE\ g.\ is-inverse_L\ f\ g)$

lemma $inv_L\text{-eq}$:

assumes $is-inverse_L\ f\ g$

shows $inv_L\ f = g$

$\langle proof \rangle$

lemma $inv_L\text{-I}$:

assumes $f\ o_L\ g = id\text{-blinfun}\ g\ o_L\ f = id\text{-blinfun}$

shows $g = inv_L\ f$

$\langle proof \rangle$

lemma $inv\text{-app1}\ [simp]:\ invertible_L\ X \implies (inv_L\ X)\ o_L\ X = id\text{-blinfun}$

$\langle proof \rangle$

lemma $inv\text{-app2}\ [simp]:\ invertible_L\ X \implies X\ o_L\ (inv_L\ X) = id\text{-blinfun}$

$\langle proof \rangle$

lemma $inv\text{-app1}'\ [simp]:\ invertible_L\ X \implies inv_L\ X\ (X\ v) = v$

$\langle proof \rangle$

lemma $inv\text{-app2}'\ [simp]:\ invertible_L\ X \implies X\ (inv_L\ X\ v) = v$

$\langle proof \rangle$

lemma $[simp]:\ invertible_L\ X \implies inv_L\ (inv_L\ X) = X$

$\langle proof \rangle$

lemma $inv_L\text{-cancel-iff}$:

assumes $invertible_L\ f$

shows $f\ x = y \iff x = inv_L\ f\ y$

$\langle proof \rangle$

lemma $invertible_L\text{-inf-sum}$:

assumes $norm\ (X :: 'b :: banach \Rightarrow_L\ 'b) < 1$

shows $invertible_L\ (id\text{-blinfun} - X)$

$\langle proof \rangle$

lemma $inv_L\text{-inf-sum}$:

fixes $X :: 'b :: banach \Rightarrow_L -$

assumes $norm\ X < 1$

shows $inv_L\ (id\text{-blinfun} - X) = (\sum\ i.\ X\ \hat{\sim}\ i)$

$\langle proof \rangle$

lemma $is-inverse_L\text{-compose}$:

assumes $invertible_L\ f\ invertible_L\ g$

shows $is-inverse_L\ (f\ o_L\ g)\ (inv_L\ g\ o_L\ inv_L\ f)$

$\langle proof \rangle$

lemma *invertible_L-compose*: $invertible_L f \implies invertible_L g \implies invertible_L (f \circ_L g)$
 ⟨proof⟩

lemma *inv_L-compose*:
assumes $invertible_L f \ invertible_L g$
shows $inv_L (f \circ_L g) = (inv_L g) \circ_L (inv_L f)$
 ⟨proof⟩

lemma *inv_L-id-blinfun[simp]*: $inv_L \ id\text{-blinfun} = id\text{-blinfun}$
 ⟨proof⟩

2.5 Norm

lemma *bounded-range-subset*:
 $bounded (range f :: real \ set) \implies bounded (f \ 'X)$
 ⟨proof⟩

lemma *bounded-const*: $bounded ((\lambda \cdot. x) \ 'X)$
 ⟨proof⟩

lift-definition *bfun-pos* :: $('d \Rightarrow_b real) \Rightarrow ('d \Rightarrow_b real)$ **is** $\lambda f \ i. \text{if } f \ i < 0 \text{ then } -f \ i \text{ else } f \ i$
 ⟨proof⟩

lemma *bfun-pos-zero[simp]*: $bfun\text{-pos } f = 0 \longleftrightarrow f = 0$
 ⟨proof⟩

lift-definition *bfun-nonneg* :: $('d \Rightarrow_b real) \Rightarrow ('d \Rightarrow_b real)$ **is** $\lambda f \ i. \text{if } f \ i \leq 0 \text{ then } 0 \text{ else } f \ i$
 ⟨proof⟩

lemma *bfun-nonneg-split*: $bfun\text{-nonneg } x - bfun\text{-nonneg } (-x) = x$
 ⟨proof⟩

lemma *blinfun-split*: $blinfun\text{-apply } f \ x = f (bfun\text{-nonneg } x) - f (bfun\text{-nonneg } (-x))$
 ⟨proof⟩

lemma *bfun-nonneg-pos*: $bfun\text{-nonneg } x + bfun\text{-nonneg } (-x) = bfun\text{-pos } x$
 ⟨proof⟩

lemma *bfun-nonneg*: $0 \leq bfun\text{-nonneg } f$
 ⟨proof⟩

lemma *bfun-pos-eq-nonneg*: $bfun\text{-pos } n = bfun\text{-nonneg } n + bfun\text{-nonneg } (-n)$
 ⟨proof⟩

lemma *blinfun-mono-norm-pos*:

fixes $f :: ('c \Rightarrow_b \text{real}) \Rightarrow_L ('d \Rightarrow_b \text{real})$
assumes $\bigwedge v :: 'c \Rightarrow_b \text{real}. v \geq 0 \implies f v \geq 0$
shows $\text{norm } (f n) \leq \text{norm } (f (\text{bfun-pos } n))$
<proof>

lemma *norm-bfun-pos[simp]*: $\text{norm } (\text{bfun-pos } f) = \text{norm } f$
<proof>

lemma *norm-blinfun-mono-eq-nonneg*:

fixes $f :: ('c \Rightarrow_b \text{real}) \Rightarrow_L ('d \Rightarrow_b \text{real})$
assumes $\bigwedge v :: 'c \Rightarrow_b \text{real}. v \geq 0 \implies f v \geq 0$
shows $\text{norm } f = (\bigsqcup v \in \{v. v \geq 0\}. \text{norm } (f v) / \text{norm } v)$
<proof>

lemma *norm-blinfun-normalized-le*: $\text{norm } (\text{blinfun-apply } f v) / \text{norm } v \leq \text{norm } f$
<proof>

lemma *norm-blinfun-mono-eq-nonneg'*:

fixes $f :: ('c \Rightarrow_b \text{real}) \Rightarrow_L ('d \Rightarrow_b \text{real})$
assumes $\bigwedge v :: 'c \Rightarrow_b \text{real}. 0 \leq v \implies 0 \leq f v$
shows $\text{norm } f = (\bigsqcup x \in \{x. \text{norm } x = 1 \wedge x \geq 0\}. \text{norm } (f x))$
<proof>

lemma *norm-blinfun-mono-le-norm-one*:

fixes $f :: ('c \Rightarrow_b \text{real}) \Rightarrow_L ('d \Rightarrow_b \text{real})$
assumes $\bigwedge v :: 'c \Rightarrow_b \text{real}. v \geq 0 \implies f v \geq 0$
assumes $\text{norm } x = 1 \implies 0 \leq x$
shows $\text{norm } (f x) \leq \text{norm } (f 1)$
<proof>

lemma *norm-blinfun-mono-eq-one*:

fixes $f :: ('c \Rightarrow_b \text{real}) \Rightarrow_L ('d \Rightarrow_b \text{real})$
assumes $\bigwedge v :: 'c \Rightarrow_b \text{real}. v \geq 0 \implies f v \geq 0$
shows $\text{norm } f = \text{norm } (f 1)$
<proof>

2.6 Miscellaneous

lemma *bounded-linear-apply-bfun*: $\text{bounded-linear } (\lambda x. \text{apply-bfun } x i)$
<proof>

lemma *lim-blinfun-apply*: $\text{convergent } X \implies (\lambda n. \text{blinfun-apply } (X n) u) \longrightarrow \text{lim } X u$
<proof>

lemma *bounded-apply-blinfun*:

```

assumes bounded ((F :: 'c  $\Rightarrow$  'd::real-normed-vector  $\Rightarrow_L$  'b::real-normed-vector)
' S)
shows bounded (( $\lambda$ b. blinfun-apply (F b) x) ' S)
<proof>
end
theory MDP-reward-Util
imports Blinfun-Util
begin

```

3 Auxiliary Lemmas

3.1 Summability

```

lemma summable-powser-const:
fixes c :: real
assumes |c| < 1
shows summable ( $\lambda$ n. cn * x)
<proof>

```

3.2 Infinite sums

```

lemma suminf-split-head':
summable (f :: nat  $\Rightarrow$  'x :: real-normed-vector)  $\implies$  suminf f = f 0
+ ( $\sum$  n. f (Suc n))
<proof>

```

```

lemma sum-disc-lim:
assumes |c :: real| < 1
shows ( $\sum$  x. cn * B) = B / (1-c)
<proof>

```

3.3 Bounded Functions

```

lemma suminf-apply-bfun:
fixes f :: nat  $\Rightarrow$  'c  $\Rightarrow_b$  real
assumes summable f
shows ( $\sum$  i. f i) x = ( $\sum$  i. f i x)
<proof>

```

```

lemma sum-apply-bfun:
fixes f :: nat  $\Rightarrow$  'c  $\Rightarrow_b$  real
shows ( $\sum$  i<n. f i) x = ( $\sum$  i<n. apply-bfun (f i) x)
<proof>

```

3.4 Push-Forward of a Bounded Function

```

lemma integrable-bfun-prob-space [simp]:
integrable (measure-pmf P) ( $\lambda$ t. apply-bfun f (F t) :: real)
<proof>

```


lift-definition *push-exp* :: ('b ⇒ 'c pmf) ⇒ ('c ⇒_b real) ⇒ ('b ⇒_b real) **is**
 $\lambda c f s. \text{measure-pmf.expectation } (c \ s) \ f$
 ⟨proof⟩

declare *push-exp.rep-eq*[simp]

lemma *norm-push-exp-le-norm*: $\text{norm } (\text{push-exp } d \ x) \leq \text{norm } x$
 ⟨proof⟩

lemma *push-exp-bounded-linear* [simp]: *bounded-linear* (*push-exp* *d*)
 ⟨proof⟩

lemma *onorm-push-exp* [simp]: *onorm* (*push-exp* *d*) = 1
 ⟨proof⟩

lemma *push-exp-return*[simp]: *push-exp* *return-pmf* = *id*
 ⟨proof⟩

3.5 Boundedness

lemma *bounded-abs*[intro]:
 $\text{bounded } (X' :: \text{real set}) \implies \text{bounded } (\text{abs } ' X')$
 ⟨proof⟩

lemma *bounded-abs-range*[intro]:
 $\text{bounded } (\text{range } f :: \text{real set}) \implies \text{bounded } (\text{range } (\lambda x. \text{abs } (f \ x)))$
 ⟨proof⟩

3.6 Probability Theory

lemma *integral-measure-pmf-bind*:
assumes $(\bigwedge x. |\text{f } x :: 'b \Rightarrow \text{real } x| \leq B)$
shows $(\int x. \text{f } x \ \partial((\text{measure-pmf } M) \gg (\lambda x. \text{measure-pmf } (N \ x))))$
 $= (\int x. \int y. \text{f } y \ \partial N \ x \ \partial M)$
 ⟨proof⟩

lemma *lemma-4-3-1'*:
assumes $\text{set-pmf } p \subseteq W$
and $\text{bounded } ((w :: 'c \Rightarrow \text{real}) ' W)$
and $W \neq \{\}$
and $\text{measure-pmf.expectation } p \ w = (\bigsqcup p \in \{p. \text{set-pmf } p \subseteq W\}.$
 $\text{measure-pmf.expectation } p \ w)$
shows $\exists x \in W. \text{measure-pmf.expectation } p \ w = w \ x$
 ⟨proof⟩

lemma *lemma-4-3-1*:
assumes $\text{set-pmf } p \subseteq W$ *integrable* (*measure-pmf* *p*) *w* *bounded* $((w :: 'c \Rightarrow \text{real}) ' W)$

shows *measure-pmf.expectation* $p \ w \leq \bigsqcup (w \text{ ' } W)$
 ⟨proof⟩

lemma *bounded-integrable*:

assumes *bounded* (*range* v) $v \in \text{borel-measurable}$ (*measure-pmf* p)
shows *integrable* (*measure-pmf* p) ($v :: \text{'c} \Rightarrow \text{real}$)
 ⟨proof⟩

3.7 Argmax

lemma *finite-is-arg-max*: *finite* $X \Rightarrow X \neq \{\}$ $\Rightarrow \exists x. \text{is-arg-max}$ ($f :: \text{'c} \Rightarrow \text{real}$) ($\lambda x. x \in X$) x
 ⟨proof⟩

lemma *finite-arg-max-le*:

assumes *finite* ($X :: \text{'c set}$) $X \neq \{\}$
shows $s \in X \Rightarrow (f :: \text{'c} \Rightarrow \text{real}) \ s \leq f$ (*arg-max-on* ($f :: \text{'c} \Rightarrow \text{real}$) X)
 ⟨proof⟩

lemma *arg-max-on-in*:

assumes *finite* ($X :: \text{'c set}$) $X \neq \{\}$
shows (*arg-max-on* ($f :: \text{'c} \Rightarrow \text{real}$) X) $\in X$
 ⟨proof⟩

lemma *finite-arg-max-eq-Max*:

assumes *finite* ($X :: \text{'c set}$) $X \neq \{\}$
shows ($f :: \text{'c} \Rightarrow \text{real}$) (*arg-max-on* f X) = *Max* ($f \text{ ' } X$)
 ⟨proof⟩

lemma *arg-max-SUP*: *is-arg-max* ($f :: \text{'b} \Rightarrow \text{real}$) ($\lambda x. x \in X$) $m \Rightarrow f \ m = (\bigsqcup (f \text{ ' } X))$
 ⟨proof⟩

definition *has-max* $X \equiv \exists x \in X. \forall x' \in X. x' \leq x$

definition *has-arg-max* f $X \equiv \exists x. \text{is-arg-max}$ f ($\lambda x. x \in X$) x

lemma *has-max* ($(f :: \text{'b} \Rightarrow \text{real}) \text{ ' } X$) $\longleftrightarrow \text{has-arg-max}$ f X
 ⟨proof⟩

lemma *has-arg-max-is-arg-max*: *has-arg-max* f $X \Rightarrow \text{is-arg-max}$ f ($\lambda x. x \in X$) (*arg-max* f ($\lambda x. x \in X$))
 ⟨proof⟩

lemma *has-arg-max-arg-max*: *has-arg-max* f $X \Rightarrow (\text{arg-max}$ f ($\lambda x. x \in X$)) $\in X$
 ⟨proof⟩

lemma *app-arg-max-ge*: *has-arg-max* ($f :: 'b \Rightarrow \text{real}$) $X \Longrightarrow x \in X$
 $\Longrightarrow f\ x \leq f\ (\text{arg-max-on } f\ X)$
 $\langle \text{proof} \rangle$

lemma *app-arg-max-eq-SUP*: *has-arg-max* ($f :: 'b \Rightarrow \text{real}$) $X \Longrightarrow f$
 $(\text{arg-max-on } f\ X) = \bigsqcup (f\ 'X)$
 $\langle \text{proof} \rangle$

lemma *SUP-is-arg-max*:
assumes $x \in X$ *bdd-above* ($f\ 'X$) ($f :: 'c \Rightarrow \text{real}$) $x = \bigsqcup (f\ 'X)$
shows *is-arg-max* $f\ (\lambda x. x \in X)\ x$
 $\langle \text{proof} \rangle$

lemma *is-arg-max-linorderI*[*intro*]: **fixes** $f :: 'c \Rightarrow 'b :: \text{linorder}$
assumes $P\ x \wedge y. (P\ y \Longrightarrow f\ x \geq f\ y)$
shows *is-arg-max* $f\ P\ x$
 $\langle \text{proof} \rangle$

lemma *is-arg-max-linorderD*[*dest*]: **fixes** $f :: 'c \Rightarrow 'b :: \text{linorder}$
assumes *is-arg-max* $f\ P\ x$
shows $P\ x (P\ y \Longrightarrow f\ x \geq f\ y)$
 $\langle \text{proof} \rangle$

lemma *is-arg-max-cong*:
assumes $\bigwedge x. P\ x \Longrightarrow f\ x = g\ x$
shows *is-arg-max* $f\ P\ x \longleftrightarrow \text{is-arg-max } g\ P\ x$
 $\langle \text{proof} \rangle$

lemma *is-arg-max-congI*:
assumes *is-arg-max* $f\ P\ x \wedge x. P\ x \Longrightarrow f\ x = g\ x$
shows *is-arg-max* $g\ P\ x$
 $\langle \text{proof} \rangle$

3.8 Contraction Mappings

definition *is-contraction* $C \equiv \exists l. 0 \leq l \wedge l < 1 \wedge (\forall v\ u. \text{dist } (C\ v)$
 $(C\ u) \leq l * \text{dist } v\ u)$

lemma *banach'*:
fixes $C :: 'b :: \text{complete-space} \Rightarrow 'b$
assumes *is-contraction* C
shows $\exists! v. C\ v = v \wedge v. (\lambda n. (C \rightsquigarrow n)\ v) \longrightarrow (\text{THE } v. C\ v = v)$
 $\langle \text{proof} \rangle$

lemma *contraction-dist*:
fixes $C :: 'b :: \text{complete-space} \Rightarrow 'b$
assumes $\bigwedge v\ u. \text{dist } (C\ v)\ (C\ u) \leq c * \text{dist } v\ u$
assumes $0 \leq c < 1$

shows $(1 - c) * \text{dist } v \text{ (THE } v. C v = v) \leq \text{dist } v \text{ (} C v)$
 ⟨proof⟩

3.9 Limits

lemma *tendsto-bfun-sandwich*:

assumes

$(f :: \text{nat} \Rightarrow 'b \Rightarrow_b \text{real}) \longrightarrow x$ $(g :: \text{nat} \Rightarrow 'b \Rightarrow_b \text{real}) \longrightarrow x$
eventually $(\lambda n. f n \leq h n)$ *sequentially eventually* $(\lambda n. h n \leq g n)$

sequentially

shows $(h :: \text{nat} \Rightarrow 'b \Rightarrow_b \text{real}) \longrightarrow x$

⟨proof⟩

3.10 Supremum

lemma *SUP-add-le*:

assumes $X \neq \{\}$ *bounded* $(B ' X)$ *bounded* $(A' ' X)$

shows $(\bigsqcup c \in X. (B :: 'a \Rightarrow \text{real}) c + A' c) \leq (\bigsqcup b \in X. B b) +$
 $(\bigsqcup a \in X. A' a)$

⟨proof⟩

lemma *le-SUP-diff'*:

assumes $ne: X \neq \{\}$

and *bdd*: *bounded* $(B ' X)$ *bounded* $(A' ' X)$

and *sup-le*: $(\bigsqcup a \in X. (A' :: 'a \Rightarrow \text{real}) a) \leq (\bigsqcup b \in X. B b)$

shows $(\bigsqcup b \in X. B b) - (\bigsqcup a \in X. (A' :: 'a \Rightarrow \text{real}) a) \leq (\bigsqcup c \in$
 $X. B c - A' c)$

⟨proof⟩

lemma *le-SUP-diff*:

fixes $A' :: 'a \Rightarrow \text{real}$

assumes $X \neq \{\}$ *bounded* $(B ' X)$ *bounded* $(A' ' X)$ $(\bigsqcup a \in X. A'$
 $a) \leq (\bigsqcup b \in X. B b)$

shows $0 \leq (\bigsqcup c \in X. B c - A' c)$

⟨proof⟩

lemma *bounded-SUP-mul[simp]*:

$X \neq \{\} \Longrightarrow 0 \leq l \Longrightarrow \text{bounded } (f ' X) \Longrightarrow (\bigsqcup x \in X. (l :: \text{real}) * f$
 $x) = (l * (\bigsqcup x \in X. f x))$

⟨proof⟩

lemma *abs-cSUP-le[intro]*:

$X \neq \{\} \Longrightarrow \text{bounded } (F ' X) \Longrightarrow |\bigsqcup x \in X. (F x) :: \text{real}| \leq (\bigsqcup x \in$
 $X. |F x|)$

⟨proof⟩

end

4 Discrete-Time Markov Decision Processes with Arbitrary State Spaces

In this file we construct discrete-time Markov decision processes, e.g. with arbitrary state spaces. Proofs and definitions are adapted from `Markov_Models.Discrete_Time_Markov_Process`.

theory *MDP-cont*

imports *HOL-Probability.Probability*

begin

lemma *Ionescu-Tulcea-C-eq*:

assumes $\bigwedge i h. h \in \text{space } (PiM \{0..<i\} N) \implies P i h = P' i h$

assumes *h*: *Ionescu-Tulcea* *P N Ionescu-Tulcea P' N*

shows *Ionescu-Tulcea.C P N 0 n* ($\lambda x. \text{undefined}$) = *Ionescu-Tulcea.C P' N 0 n* ($\lambda x. \text{undefined}$)

<proof>

lemma *Ionescu-Tulcea-CI-eq*:

assumes $\bigwedge i h. h \in \text{space } (PiM \{0..<i\} N) \implies P i h = P' i h$

assumes *h*: *Ionescu-Tulcea P N Ionescu-Tulcea P' N*

shows *Ionescu-Tulcea.CI P N* = *Ionescu-Tulcea.CI P' N*

<proof>

lemma *measure-eqI-PiM-sequence*:

fixes *M* :: *nat* \Rightarrow *'a measure*

assumes **[simp]*: *sets P* = *PiM UNIV M sets Q* = *PiM UNIV M*

assumes *eq*: $\bigwedge A n. (\bigwedge i. A i \in \text{sets } (M i)) \implies$

$P (\text{prod-emb UNIV M } \{..n\} (Pi_E \{..n\} A)) = Q (\text{prod-emb UNIV M } \{..n\} (Pi_E \{..n\} A))$

assumes *A*: *finite-measure P*

shows *P* = *Q*

<proof>

lemma *distr-cong-simp*:

$M = K \implies \text{sets } N = \text{sets } L \implies (\bigwedge x. x \in \text{space } M = \text{simp} \Rightarrow f x = g x) \implies \text{distr } M N f = \text{distr } K L g$

<proof>

4.1 Definition and Basic Properties

locale *discrete-MDP* =

fixes *Ms* :: *'s measure*

and *Ma* :: *'a measure*

and *A* :: *'s* \Rightarrow *'a set*

and *K* :: *'s* \times *'a* \Rightarrow *'s measure*

assumes *A-s*: $\bigwedge s. A s \in \text{sets } Ma$

assumes $A\text{-ne}$: $\bigwedge s. A\ s \neq \{\}$

assumes $ex\text{-pol}$: $\exists \delta \in Ms \rightarrow_M Ma. \forall s. \delta\ s \in A\ s$

assumes $K[\text{measurable}]$: $K \in Ms \otimes_M Ma \rightarrow_M \text{prob-algebra } Ms$
begin

lemma $space\text{-prodI}[\text{intro}]$: $x \in \text{space } A' \implies y \in \text{space } B \implies (x,y) \in \text{space } (A' \otimes_M B)$
 $\langle \text{proof} \rangle$

abbreviation $M \equiv Ms \otimes_M Ma$

abbreviation $Ma\text{-}A\ s \equiv \text{restrict-space } Ma\ (A\ s)$

lemma $space\text{-ma}[\text{intro}]$: $s \in \text{space } Ms \implies a \in \text{space } Ma \implies (s,a) \in \text{space } M$
 $\langle \text{proof} \rangle$

lemma $space\text{-x0}[\text{simp}]$: $x0 \in \text{space } (\text{prob-algebra } Ms) \implies \text{space } x0 = \text{space } Ms$
 $\langle \text{proof} \rangle$

lemma $A\text{-subs-}Ma$: $A\ s \subseteq \text{space } Ma$
 $\langle \text{proof} \rangle$

lemma $space\text{-}Ma\text{-}A\text{-subset}$: $s \in \text{space } Ms \implies \text{space } (Ma\text{-}A\ s) \subseteq A\ s$
 $\langle \text{proof} \rangle$

lemma $K\text{-restrict } [\text{measurable}]$: $K \in (Ms \otimes_M Ma\text{-}A\ s) \rightarrow_M \text{prob-algebra } Ms$
 $\langle \text{proof} \rangle$

lemma $measurable\text{-}K\text{-act}[\text{measurable}, \text{intro}]$: $s \in \text{space } Ms \implies (\lambda a. K\ (s, a)) \in Ma \rightarrow_M \text{prob-algebra } Ms$
 $\langle \text{proof} \rangle$

lemma $measurable\text{-}K\text{-st}[\text{measurable}, \text{intro}]$: $a \in \text{space } Ma \implies (\lambda s. K\ (s, a)) \in Ms \rightarrow_M \text{prob-algebra } Ms$
 $\langle \text{proof} \rangle$

lemma $space\text{-}K[\text{simp}]$: $sa \in \text{space } M \implies \text{space } (K\ sa) = \text{space } Ms$
 $\langle \text{proof} \rangle$

lemma $space\text{-}K2[\text{simp}]$: $s \in \text{space } Ms \implies a \in \text{space } Ma \implies \text{space } (K\ (s, a)) = \text{space } Ms$
 $\langle \text{proof} \rangle$

lemma $space\text{-}K\text{-}E$: $s' \in \text{space } (K\ (s,a)) \implies s \in \text{space } Ms \implies a \in \text{space } Ma \implies s' \in \text{space } Ms$

$\langle proof \rangle$

lemma *sets-K*: $sa \in space\ M \implies sets\ (K\ sa) = sets\ Ms$
 $\langle proof \rangle$

lemma *sets-K'*[*measurable-cong*]: $s \in space\ Ms \implies a \in space\ Ma \implies$
 $sets\ (K\ (s,a)) = sets\ Ms$
 $\langle proof \rangle$

lemma *prob-space-K*[*intro*]: $sa \in space\ M \implies prob-space\ (K\ sa)$
 $\langle proof \rangle$

lemma *prob-space-K2*[*intro*]: $s \in space\ Ms \implies a \in space\ Ma \implies$
 $prob-space\ (K\ (s,a))$
 $\langle proof \rangle$

lemma *K-in-space* [*intro*]: $m \in space\ M \implies K\ m \in space\ (prob-algebra\ Ms)$
 $\langle proof \rangle$

4.2 Policies

type-synonym (*'c, 'd*) *pol* = $nat \Rightarrow ((nat \Rightarrow 'c \times 'd) \times 'c) \Rightarrow 'd$
measure

abbreviation $H\ i \equiv Pi_M\ \{0..<i\}\ (\lambda-. M)$

abbreviation $Hs\ i \equiv H\ i \otimes_M Ms$

lemma *space-H1*: $j < (i :: nat) \implies \omega \in space\ (H\ i) \implies \omega\ j \in space\ M$
 $\langle proof \rangle$

lemma *space-case-nat*[*intro*]:
assumes $\omega \in space\ (H\ i)\ s \in space\ Ms$
shows *case-nat* $s\ (fst \circ \omega)\ i \in space\ Ms$
 $\langle proof \rangle$

lemma *undefined-in-H0*: $(\lambda-. undefined) \in space\ (H\ (0 :: nat))$
 $\langle proof \rangle$

lemma *sets-K-Suc*[*measurable-cong*]: $h \in space\ (H\ (Suc\ n)) \implies sets\ (K\ (h\ n)) = sets\ Ms$
 $\langle proof \rangle$

A decision rule is a function from states to distributions over enabled actions.

definition *is-dec0* $d \equiv d \in Ms \rightarrow_M prob-algebra\ Ma$

definition *is-dec* ($t :: \text{nat}$) $d \equiv (d \in \text{Hs } t \rightarrow_M \text{prob-algebra } Ma)$

lemma *is-dec0* $d \implies \text{is-dec } t (\lambda(-,s). d s)$
 $\langle \text{proof} \rangle$

A policy is a function from histories to valid decision rules.

definition *is-policy* $:: ('s, 'a) \text{pol} \Rightarrow \text{bool}$ **where**
is-policy $p \equiv \forall i. \text{is-dec } i (p i)$

abbreviation *p0* $:: ('s, 'a) \text{pol} \Rightarrow 's \Rightarrow 'a \text{ measure}$ **where**
p0 $p s \equiv p (0 :: \text{nat}) (\lambda-. \text{undefined}, s)$

context

fixes p **assumes** $p[\text{simp}]$: *is-policy* p
begin

lemma *is-policyD[measurable]*: $p i \in \text{Hs } i \rightarrow_M \text{prob-algebra } Ma$
 $\langle \text{proof} \rangle$

lemma *space-policy[simp]*: $hs \in \text{space } (Hs i) \implies \text{space } (p i hs) = \text{space } Ma$
 $\langle \text{proof} \rangle$

lemma *space-policy'[simp]*: $h \in \text{space } (H i) \implies s \in \text{space } Ms \implies \text{space } (p i (h,s)) = \text{space } Ma$
 $\langle \text{proof} \rangle$

lemma *space-policyI[intro]*:
assumes $s \in \text{space } Ms$ $h \in \text{space } (H i)$ $a \in \text{space } Ma$
shows $a \in \text{space } (p i (h,s))$
 $\langle \text{proof} \rangle$

lemma *sets-policy[simp]*: $hs \in \text{space } (Hs i) \implies \text{sets } (p i hs) = \text{sets } Ma$
 $\langle \text{proof} \rangle$

lemma *sets-policy'[measurable-cong, simp]*:
 $h \in \text{space } (H i) \implies s \in \text{space } Ms \implies \text{sets } (p i (h,s)) = \text{sets } Ma$
 $\langle \text{proof} \rangle$

lemma *sets-policy''[measurable-cong, simp]*:
 $h \in \text{space } ((Pi_M \{ \} (\lambda-. M))) \implies s \in \text{space } Ms \implies \text{sets } (p 0 (h,s)) = \text{sets } Ma$
 $\langle \text{proof} \rangle$

lemma *policy-prob-space*: $hs \in \text{space } (Hs i) \implies \text{prob-space } (p i hs)$
 $\langle \text{proof} \rangle$

lemma *policy-prob-space'*: $h \in \text{space } (H \ i) \implies s \in \text{space } Ms \implies \text{prob-space } (p \ i \ (h,s))$
 ⟨proof⟩

lemma *prob-space-p0*: $x \in \text{space } Ms \implies \text{prob-space } (p0 \ p \ x)$
 ⟨proof⟩

lemma *p0-sets[measurable-cong]*: $x \in \text{space } Ms \implies \text{sets } (p \ 0 \ (\lambda \cdot \text{undefined},x)) = \text{sets } Ma$
 ⟨proof⟩

lemma *space-p0[simp]*: $s \in \text{space } Ms \implies \text{space } (p0 \ p \ s) = \text{space } Ma$
 ⟨proof⟩

lemma *return-policy-prob-algebra [measurable]*:
 $h \in \text{space } (H \ n) \implies x \in \text{space } Ms \implies (\lambda a. \text{return } M \ (x, a)) \in p \ n$
 $(h, x) \rightarrow_M \text{prob-algebra } M$
 ⟨proof⟩
end

4.3 Successor Policy

To shift the policy by one step, we provide a single state-action pair as history

definition *Suc-policy* $p \ sa = (\lambda i \ (h, s). p \ (Suc \ i) \ (\lambda i'. \text{case-nat } sa \ h \ i', s))$

lemma *p-as-Suc-policy*: $p \ (Suc \ i) \ (h, s) = \text{Suc-policy } p \ ((h \ 0)) \ i \ (\lambda i. h \ (Suc \ i), s)$
 ⟨proof⟩

lemma *is-policy-Suc-policy[intro]*:
assumes $s: sa \in \text{space } M$ **and** $p: \text{is-policy } p$
shows $\text{is-policy } (\text{Suc-policy } p \ sa)$
 ⟨proof⟩

lemma *Suc-policy-measurable-step[measurable]*:
assumes $\text{is-policy } p$
shows $(\lambda x. \text{Suc-policy } p \ (fst \ (fst \ x)) \ n \ (snd \ (fst \ x), \ snd \ x)) \in$
 $(M \otimes_M Pi_M \ \{0..<n\} \ (\lambda \cdot. M)) \otimes_M Ms \rightarrow_M \text{prob-algebra } Ma$
 ⟨proof⟩

4.4 Single-Step Distribution

K' takes a policy, a distribution over 's, the epoch, and a history, produces a distribution over the next state-action pair.

definition $K' :: ('s, 'a) \text{pol} \Rightarrow 's \text{measure} \Rightarrow \text{nat} \Rightarrow (\text{nat} \Rightarrow ('s \times 'a))$

$\Rightarrow ('s \times 'a)$ *measure*
where
 $K' p s0 n \omega = do \{$
 $s \leftarrow case\text{-}nat\ s0\ (K \circ \omega)\ n;$
 $a \leftarrow p\ n\ (\omega, s);$
 $return\ M\ (s, a)$
 $\}$

lemma *prob-space-K'*:
assumes p : *is-policy* p **and** x : $x0 \in space\ (prob\text{-}algebra\ Ms)$ **and** h :
 $h \in space\ (H\ n)$
shows *prob-space* $(K' p x0 n h)$
 $\langle proof \rangle$

lemma *measurable-K'[measurable]*:
assumes p : *is-policy* p **and** x : $x \in space\ (prob\text{-}algebra\ Ms)$
shows $K' p x i \in H\ i \rightarrow_M\ prob\text{-}algebra\ M$
 $\langle proof \rangle$

4.5 Initial State-Action Distribution

$K0$ produces the initial state-action distribution from a state distribution and a policy.

definition $K0\ p\ s0 = K' p\ s0\ 0\ (\lambda\ \cdot.\ undefined)$

lemma *K0-def'*:
 $K0\ p\ s0 = do \{$
 $s \leftarrow s0;$
 $a \leftarrow p0\ p\ s;$
 $return\ M\ (s, a)\}$
 $\langle proof \rangle$

lemma *K0-prob[measurable]*: *is-policy* $p \implies K0\ p \in prob\text{-}algebra\ Ms$
 $\rightarrow_M\ prob\text{-}algebra\ M$
 $\langle proof \rangle$

lemma *prob-space-K0*: *is-policy* $p \implies x0 \in space\ (prob\text{-}algebra\ Ms)$
 $\implies prob\text{-}space\ (K0\ p\ x0)$
 $\langle proof \rangle$

lemma *space-K0[simp]*: *is-policy* $p \implies s \in space\ (prob\text{-}algebra\ Ms)$
 $\implies space\ (K0\ p\ s) = space\ M$
 $\langle proof \rangle$

lemma *sets-K0[measurable-cong]*:
assumes *is-policy* $p\ s \in space\ (prob\text{-}algebra\ Ms)$
shows *sets* $(K0\ p\ s) = sets\ M$
 $\langle proof \rangle$

lemma *K0-return-eq-p0*:

assumes *is-policy* p $s \in \text{space } Ms$

shows $K0\ p\ (\text{return } Ms\ s) = p0\ p\ s \gg\! \gg (\lambda a. \text{return } M\ (s,a))$

$\langle \text{proof} \rangle$

lemma *M-ne-policy[intro]*: *is-policy* $p \implies s \in \text{space } (\text{prob-algebra } Ms)$

$\implies \text{space } M \neq \{\}$

$\langle \text{proof} \rangle$

4.6 Sequence Space of the MDP

We can instantiate *Ionescu-Tulcea* with K' .

lemma *IT-K'*: *is-policy* $p \implies x \in \text{space } (\text{prob-algebra } Ms) \implies \text{Ionescu-Tulcea}$

$(K'\ p\ x)\ (\lambda-. M)$

$\langle \text{proof} \rangle$

definition *lim-sequence* :: $('s, 'a)\ \text{pol} \Rightarrow 's\ \text{measure} \Rightarrow (\text{nat} \Rightarrow ('s \times 'a))\ \text{measure}$

where

lim-sequence $p\ x = \text{projective-family.lim } UNIV\ (\text{Ionescu-Tulcea.CI}\ (K'\ p\ x)\ (\lambda-. M))\ (\lambda-. M)$

lemma

assumes $x: x \in \text{space } (\text{prob-algebra } Ms)$ **and** $p: \text{is-policy } p$

shows *space-lim-sequence*: $\text{space } (\text{lim-sequence } p\ x) = \text{space } (\Pi_M\ i \in UNIV. M)$

and *sets-lim-sequence[measurable-cong]*: $\text{sets } (\text{lim-sequence } p\ x) = \text{sets } (\Pi_M\ i \in UNIV. M)$

and *emeasure-lim-sequence-emb*: $\bigwedge J\ X. \text{finite } J \implies X \in \text{sets } (\Pi_M\ j \in J. M) \implies$

$\text{emeasure } (\text{lim-sequence } p\ x)\ (\text{prod-emb } UNIV\ (\lambda-. M)\ J\ X) = \text{emeasure } (\text{Ionescu-Tulcea.CI}\ (K'\ p\ x)\ (\lambda-. M)\ J)\ X$

and *emeasure-lim-sequence-emb-I0o*: $\bigwedge n\ X. X \in \text{sets } (\Pi_M\ i \in \{0..<n\}. M) \implies$

$\text{emeasure } (\text{lim-sequence } p\ x)\ (\text{prod-emb } UNIV\ (\lambda-. M)\ \{0..<n\}\ X) =$

$\text{emeasure } (\text{Ionescu-Tulcea.C}\ (K'\ p\ x)\ (\lambda-. M)\ 0\ n\ (\lambda x. \text{undefined}))\ X$

$\langle \text{proof} \rangle$

lemma *lim-sequence-prob-space*:

assumes *is-policy* p $s \in \text{space } (\text{prob-algebra } Ms)$

shows *prob-space* $(\text{lim-sequence } p\ s)$

$\langle \text{proof} \rangle$

4.7 Measurability of the Sequence Space

lemma *lim-sequence[measurable]*:

assumes $p: \text{is-policy } p$

shows $\lim\text{-sequence } p \in \text{prob-algebra } Ms \rightarrow_M \text{prob-algebra } (\Pi_M$
 $i \in UNIV. M)$
 $\langle \text{proof} \rangle$

lemma $\lim\text{-sequence-aux}[\text{measurable}]$:

assumes p : $\text{is-policy } p$

assumes f : $\bigwedge x. x \in \text{space } M \implies \text{is-policy } (f x)$

assumes f' : $\bigwedge n. (\lambda x. f (fst (fst x)) n (snd (fst x), snd x)) \in$

$(M \otimes_M Pi_M \{0..<n\} (\lambda-. M)) \otimes_M Ms \rightarrow_M \text{prob-algebra } Ma$

assumes gm : $g \in M \rightarrow_M \text{prob-algebra } Ms$

shows $(\lambda x. \lim\text{-sequence } (f x) (g x)) \in M \rightarrow_M \text{prob-algebra } (Pi_M$
 $UNIV (\lambda-. M))$

$\langle \text{proof} \rangle$

lemma $\lim\text{-sequence-Suc-return}[\text{measurable}]$:

assumes p : $\text{is-policy } p$

assumes s : $s \in \text{space } Ms$

shows $(\lambda x. \lim\text{-sequence } (\text{Suc-policy } p (s, snd x)) (\text{return } Ms (fst$
 $x))) \in$

$M \rightarrow_M \text{prob-algebra } (Pi_M UNIV (\lambda-. M))$

$\langle \text{proof} \rangle$

lemma $\lim\text{-sequence-Suc-K}[\text{measurable}]$:

assumes $\text{is-policy } p$

shows $(\lambda x. \lim\text{-sequence } (\text{Suc-policy } p x) (K x)) \in M \rightarrow_M \text{prob-algebra}$
 $(Pi_M UNIV (\lambda-. M))$

$\langle \text{proof} \rangle$

4.8 Iteration Rule

lemma step-C :

assumes x : $x \in \text{space } (\text{prob-algebra } Ms)$ **and** p : $\text{is-policy } p$

shows $\text{Ionescu-Tulcea.C } (K' p x) (\lambda-. M) 0 1 (\lambda-. \text{undefined}) \gg=$

$\text{Ionescu-Tulcea.C } (K' p x) (\lambda-. M) 1 n =$

$K0 p x \gg= (\lambda a. \text{Ionescu-Tulcea.C } (K' p x) (\lambda-. M) 1 n (\text{case-nat}$
 $a (\lambda-. \text{undefined})))$

$\langle \text{proof} \rangle$

lemma $\lim\text{-sequence-eq}$:

assumes x : $x \in \text{space } (\text{prob-algebra } Ms)$ **assumes** p : $\text{is-policy } p$

shows $\lim\text{-sequence } p x =$

$K0 p x \gg= (\lambda y. \text{distr } (\lim\text{-sequence } (\text{Suc-policy } p y) (K y)) (\Pi_M$
 $- \in UNIV. M) (\text{case-nat } y))$

$(\text{is } - = ?B p x)$

$\langle \text{proof} \rangle$

4.9 Stream Space of the MDP

definition $\lim\text{-stream} :: ('s, 'a) \text{pol} \Rightarrow 's \text{measure} \Rightarrow ('s \times 'a) \text{stream}$
 measure

where

$\text{lim-stream } p \ x = \text{distr } (\text{lim-sequence } p \ x) \ (\text{stream-space } M) \ \text{to-stream}$

lemma *space-lim-stream*: $\text{space } (\text{lim-stream } p \ x) = \text{streams } (\text{space } M)$
<proof>

lemma *sets-lim-stream[measurable-cong]*: $\text{sets } (\text{lim-stream } p \ x) = \text{sets}$
(stream-space M)
<proof>

lemma *lim-stream[measurable]*:
assumes *is-policy p*
shows $\text{lim-stream } p \in \text{prob-algebra } Ms \rightarrow_M \text{prob-algebra } (\text{stream-space } M)$
<proof>

lemma *lim-stream-Suc[measurable]*:
assumes *p: is-policy p*
shows $(\lambda a. \text{lim-stream } (\text{Suc-policy } p \ a) \ (K \ a)) \in M \rightarrow_M \text{prob-algebra}$
(stream-space M)
<proof>

lemma *space-stream-space-M-ne*: $x \in \text{space } M \implies \text{space } (\text{stream-space } M) \neq \{\}$
<proof>

lemma *prob-space-lim-stream[intro]*:
assumes *is-policy p* $x \in \text{space } (\text{prob-algebra } Ms)$
shows $\text{prob-space } (\text{lim-stream } p \ x)$
<proof>

lemma *prob-space-step*:
assumes *is-policy p* $x \in \text{space } M$
shows $\text{prob-space } (\text{lim-stream } (\text{Suc-policy } p \ x) \ (K \ x))$
<proof>

lemma *lim-stream-eq*:
assumes *p: is-policy p*
assumes *x: x ∈ space (prob-algebra Ms)*
shows $\text{lim-stream } p \ x = \text{do } \{$
 $y \leftarrow K0 \ p \ x;$
 $\omega \leftarrow \text{lim-stream } (\text{Suc-policy } p \ y) \ (K \ y);$
 $\text{return } (\text{stream-space } M) \ (y \ \#\# \ \omega)$
 $\}$
<proof>

end
end

```

theory MDP-disc
  imports
    MDP-cont
    HOL-Library.Omega-Words-Fun
begin

```

5 Markov Decision Processes with Discrete State Spaces

```

lemma (in prob-space) integral-stream-space:
  fixes  $f :: 'a \text{ stream} \Rightarrow ('b :: \{\text{banach, second-countable-topology, real-normed-vector}\})$ 
  assumes int-f: integrable (stream-space M) f
  assumes [measurable]:  $f \in \text{borel-measurable (stream-space M)}$ 
  shows  $(\int X. f X \ \partial \text{stream-space } M) = (\int x. (\int X. f (x \#\# X) \ \partial \text{stream-space } M) \ \partial M)$ 
  <proof>

```

```

lemma prefix-cons:
  Omega-Words-Fun.prefix (Suc n) seq = seq 0 \# Omega-Words-Fun.prefix n (\lambda n. seq (Suc n))
  <proof>

```

```

lemma restrict-Suc: restrict y \{0..<Suc i\} (Suc n) = (restrict (\lambda n. y (Suc n)) \{0..<i\}) n
  <proof>

```

```

lemma prefix-restrict: Omega-Words-Fun.prefix i (restrict y \{0..<i\}) = Omega-Words-Fun.prefix i y
  <proof>

```

```

lemma prefix-measurable[measurable]:
  Omega-Words-Fun.prefix i \in Pi_M \{0..<i\} (\lambda-. count-space (UNIV :: ('s :: countable \times 'a :: countable) set)) \to_M count-space UNIV
  <proof>

```

```

no-notation Omega-Words-Fun.build (infixr <\#\#> 65)

```

```

locale discrete-MDP =
  fixes  $A :: 's :: \text{countable} \Rightarrow 'a :: \text{countable set}$  — enabled actions
  and  $K :: 's \times 'a \Rightarrow 's \text{ pmf}$  — MDP kernel, transition probabilities
  assumes
     $A\text{-ne}: \bigwedge s. A\ s \neq \{\}$  — set of enabled actions is nonempty
begin

```

5.1 Policies

Type synonym for decision rules.

type-synonym $('c, 'd) \text{ dec} = 'c \Rightarrow 'd \text{ pmf}$

definition $\text{is-dec} :: ('s, 'a) \text{ dec} \Rightarrow \text{bool}$ **where**
 $\text{is-dec } d \equiv \forall s. d \ s \subseteq A \ s$

lemma $\text{is-decI}[\text{intro}]$:
 $(\bigwedge s. \text{set-pmf } (d \ s) \subseteq A \ s) \Longrightarrow \text{is-dec } d$
 $\langle \text{proof} \rangle$

abbreviation $D_R \equiv \{d. \text{is-dec } d\}$

definition $\text{is-dec-det} :: ('s \Rightarrow 'a) \Rightarrow \text{bool}$ **where**
 $\text{is-dec-det } d \equiv \forall s. d \ s \in A \ s$

abbreviation $D_D \equiv \{d. \text{is-dec-det } d\}$

definition $\text{mk-dec-det } d \ s = \text{return-pmf } (d \ s)$

lemma $\text{is-dec-mk-dec-det-iff} [\text{simp}]$: $\text{is-dec } (\text{mk-dec-det } d) \longleftrightarrow \text{is-dec-det } d$
 $\langle \text{proof} \rangle$

lemma $D\text{-det-to-MR}[\text{intro}]$: $\text{is-dec-det } d \Longrightarrow \text{is-dec } (\text{mk-dec-det } d)$
 $\langle \text{proof} \rangle$

Due to the assumption $A \ ?s \neq \{\}$, a deterministic decision rule always exists. It immediately follows via $\text{is-dec } (\text{mk-dec-det } ?d) = \text{is-dec-det } ?d$ that a randomized decision rule also exists.

lemma SOME-is-dec-det : $\text{is-dec-det } (\lambda s. \text{SOME } a. a \in A \ s)$
 $\langle \text{proof} \rangle$

lemma $\text{ex-dec-det} [\text{simp}]$: $\exists d. \text{is-dec-det } d$
 $\langle \text{proof} \rangle$

lemma $D\text{-det-ne} [\text{simp}]$: $D_D \neq \{\}$
 $\langle \text{proof} \rangle$

lemma $D_R\text{-ne} [\text{simp}]$: $D_R \neq \{\}$
 $\langle \text{proof} \rangle$

lemma $\text{ex-dec}[\text{intro}, \text{simp}]$: $\exists d. \text{is-dec } d$
 $\langle \text{proof} \rangle$

Type synonym for policies.

type-synonym $('c, 'd) \text{ pol} = ('c \times 'd) \text{ list} \Rightarrow ('c, 'd) \text{ dec}$

A policy assigns a decision rule to each observed past.

definition *is-policy* :: ('s, 'a) pol \Rightarrow bool **where**
is-policy p $\equiv \forall hs. is_dec (p hs)$

abbreviation $\Pi_{HR} \equiv \{p. is_policy\ p\}$

Deterministic policies

definition *is-deterministic* p $\equiv is_policy\ p \wedge (\forall h\ s. \exists a. p\ h\ s = return_pmf\ a)$

definition *mk-det* p h s $\equiv return_pmf (p\ h\ s)$

abbreviation $\Pi_{HD} \equiv \{p. \forall h. p\ h \in D_D\}$

Markovian policies

definition *is-markovian* p $\equiv is_policy\ p \wedge (\forall h\ h'. length\ h = length\ h' \longrightarrow p\ h = p\ h')$

definition *mk-markovian* :: (nat \Rightarrow ('s, 'a) dec) \Rightarrow ('s, 'a) pol **where**
mk-markovian p $\equiv (\lambda h. p (length\ h))$

lemma *is-markovian-mk-iff[simp]*: *is-markovian* (mk-markovian p) \longleftrightarrow
 $(\forall n. is_dec (p\ n))$
 ⟨proof⟩

lemma *is-markovian-mk[intro]*: $\forall n. is_dec (p\ n) \Longrightarrow is_markovian$
 (mk-markovian p)
 ⟨proof⟩

lemma *mk-markovian-nil [simp]*: *mk-markovian* p [] = p 0
 ⟨proof⟩

definition *mk-markovian-det* p $\equiv (\lambda h\ s. return_pmf (p (length\ h)\ s))$

abbreviation $\Pi_{MD} \equiv \{p. \forall n. p\ n \in D_D\}$

abbreviation $\Pi_{MR} \equiv \{p. \forall n. p\ n \in D_R\}$

lemma Π_{MR} -imp-policies[*intro*]: $p \in \Pi_{MR} \Longrightarrow mk_markovian\ p \in \Pi_{HR}$
 ⟨proof⟩

lemma Π_{MD} -MR-iff[*simp*]: $(\lambda n. mk_dec_det (p\ n)) \in \Pi_{MR} \longleftrightarrow p \in \Pi_{MD}$
 ⟨proof⟩

lemma Π_{MD} -to-MR[*intro*]: $p \in \Pi_{MD} \Longrightarrow (\lambda n. mk_dec_det (p\ n)) \in \Pi_{MR}$
 ⟨proof⟩

lemma $\Pi_{MD-ne}[simp]: \Pi_{MD} \neq \{\}$
 $\langle proof \rangle$

lemma $\Pi_{MR-ne}[simp]: \Pi_{MR} \neq \{\}$
 $\langle proof \rangle$

lemma $policies-ne[simp, intro]: \Pi_{HR} \neq \{\}$
 $\langle proof \rangle$

Stationary policies

definition $is-stationary\ p \equiv is-policy\ p \wedge (\forall h\ h'.\ p\ h = p\ h')$

lemma $is-stationary-const-iff[simp]: is-stationary\ (\lambda-. d) = is-dec\ d$
 $\langle proof \rangle$

lemma $is-stationary-const[intro]: is-dec\ d \implies is-stationary\ (\lambda-. d)$
 $\langle proof \rangle$

abbreviation $mk-stationary\ p \equiv mk-markovian\ (\lambda-. p)$

abbreviation $mk-stationary-det\ d \equiv mk-markovian\ (\lambda-. mk-dec-det\ d)$

5.1.1 Successor Policy

After taking the first step in the MDP, we will know which state and which action got selected during the initial epoch. To obtain a policy that acts as if the current epoch was the initial one, we prepend the observed state-action pair to the history. The result is again a policy, i.e. it satisfies *is-policy*.

definition $\pi-Suc\ p\ sa\ h = p\ (sa\#h)$

lemma $is-policy-\pi-Suc\ [intro]: is-policy\ p \implies is-policy\ (\pi-Suc\ p\ sa)$
 $\langle proof \rangle$

lemma $Suc-mk-markovian[simp]: \pi-Suc\ (mk-markovian\ p)\ x = mk-markovian\ (\lambda n. p\ (Suc\ n))$
 $\langle proof \rangle$

5.2 Stream Space of the MDP

5.2.1 Initial State-Action Distribution

If we fix a decision rule d and an initial distribution of states $S0$, we obtain a distribution over state-action pairs in the following way: First, the initial state s is sampled from $S0$, then an action a is selected from $d\ s$.

definition $K0\ d\ S0 = do\ \{$

```

s ← S0;
a ← d s;
return-pmf (s,a)
}

```

notation $K0$ (K_0)

lemma $K0$ -iff: $K0\ d\ S0 = S0 \gg= (\lambda s. \text{map-pmf } (\lambda a. (s,a))\ (d\ s))$
 $\langle \text{proof} \rangle$

lemma $\text{vimage-pair}[simp]$: $\text{Pair } x - \{p\} = (\text{if } x = \text{fst } p \text{ then } \{\text{snd } p\}$
 $\text{else } \{\})$
 $\langle \text{proof} \rangle$

lemma $\text{pmf-K0}[simp]$: $\text{pmf } (K0\ d\ S0)\ (s,a) = \text{pmf } S0\ s * \text{pmf } (d\ s)$
 a
 $\langle \text{proof} \rangle$

lemma set-pmf-K0 : $\text{set-pmf } (K0\ p\ S0) = \{(s,a). s \in S0 \wedge a \in p\ s\}$
 $\langle \text{proof} \rangle$

lemma $\text{fst-K0}[simp]$: $\text{map-pmf } \text{fst } (K0\ p\ S0) = S0$
 $\langle \text{proof} \rangle$

abbreviation $S \equiv \text{stream-space } (\text{count-space } UNIV)$

We inherit the trace space from MDPs with continuous state-action spaces

interpretation $MDP\text{-cont}$: $MDP\text{-cont. discrete-MDP } \text{count-space } UNIV$
 $\text{count-space } UNIV\ A\ K$
 $\langle \text{proof} \rangle$

lemma $\text{count-space-M}[simp]$: $MDP\text{-cont. } M = \text{count-space } UNIV$
 $\langle \text{proof} \rangle$

lemma $\text{space-M}[simp]$: $\text{space } MDP\text{-cont. } M = UNIV$
 $\langle \text{proof} \rangle$

We reuse the stream space provided by $MDP\text{-cont. } \text{lim-stream}$

definition $T :: ('s, 'a)\ \text{pol} \Rightarrow 's\ \text{pmf} \Rightarrow ('s \times 'a)\ \text{stream measure}$
where $T\ p = MDP\text{-cont. } \text{lim-stream } (\lambda n\ (h,s).\ p\ (\text{Omega-Words-Fun. prefix } n\ h)\ s)$

lemma $\text{sets-T}[measurable-cong]$:
 $\text{sets } (T\ p\ x) = \text{sets } S$
 $\langle \text{proof} \rangle$

lemma $\text{space-stream-space-ne}[simp]$: $\text{space } S \neq \{\}$
 $\langle \text{proof} \rangle$

lemma *space-T[simp]*: $\text{space } (T \ p \ S0) = \text{space } S$
 ⟨proof⟩

lemma *is-policy-MDP-cont[intro]*:
fixes $p :: ('s \times 'a) \text{ list} \Rightarrow 's \Rightarrow 'a \text{ pmf}$
shows $\text{MDP-cont.is-policy } (\lambda n \ (h,s). \ p \ (\text{Omega-Words-Fun.prefix } n \ h) \ s)$
 ⟨proof⟩

lemma *prob-space-T[intro, simp]*: $\text{prob-space } (T \ p \ x)$
 ⟨proof⟩

lemma *T-subprob[simp]*:
 $T \ p \ S0 \in \text{space } (\text{subprob-algebra } S)$
 ⟨proof⟩

lemma *T-subprob-space [simp]*: $\text{subprob-space } (T \ p \ S0)$
 ⟨proof⟩

lemma *K0-MDP-cont-eq*:
 $\text{MDP-cont.K0 } (\lambda x \ (h,s). \ \text{measure-pmf } (p \ (\text{Omega-Words-Fun.prefix } x \ h) \ s)) \ (\text{measure-pmf } S0) =$
 $\text{K0 } (p \ []) \ S0$
 ⟨proof⟩

5.2.2 Decomposition of the Stream Space

The distribution of traces/walks the MDP allows should intuitively satisfy the following rule:

1. select the initial state s from $S0$
 2. pass it to the decision rule $p \ []$ to determine a distribution over actions
 3. select the action a
- finally pass the state-action pair (s, a) to the kernel K to get a new distribution over states $s0'$

Then the iteration repeats with the updated policy $\pi\text{-Suc } p \ (s, a)$.

The result carries over from $\llbracket \text{MDP-cont.is-policy } ?p; ?x \in \text{space } (\text{prob-algebra } (\text{count-space } \text{UNIV})) \rrbracket \Longrightarrow \text{MDP-cont.lim-stream } ?p \ ?x = \text{MDP-cont.K0 } ?p \ ?x \ggg (\lambda y. \ \text{MDP-cont.lim-stream } (\text{MDP-cont.Suc-policy } ?p \ y) \ (\text{measure-pmf } (K \ y)) \ggg (\lambda \omega. \ \text{return } (\text{stream-space } \text{MDP-cont.M}) \ (y \ \#\# \ \omega)))$.

lemma *T-eq*:

shows $T\ p\ S0 = do \{$
 $sa \leftarrow measure\text{-}pmf\ (K0\ (p\ [])\ S0);$
 $\omega \leftarrow T\ (\pi\text{-}Suc\ p\ sa)\ (K\ sa);$
 $return\ S\ (sa\ \#\#\ \omega)$
 $\}$
 $\langle proof \rangle$

lemma *T-eq-distr*:

shows $T\ p\ S0 = measure\text{-}pmf\ (K0\ (p\ [])\ S0) \gg (\lambda sa.\ distr\ (T\ (\pi\text{-}Suc\ p\ sa)\ (K\ sa))\ S\ ((\#\#\)\ sa))$
 $\langle proof \rangle$

The iteration rule lets us nicely decompose integrals (expected values) over functions on traces of the MDP.

lemma *integral-T*:

fixes $f :: ('s \times 'a)\ stream \Rightarrow real$
assumes $f\text{-bounded}$: $\bigwedge x.\ |f\ x| \leq B$
assumes f : $f \in borel\text{-}measurable\ S$
shows $(\int t.\ f\ t\ \partial T\ p\ x) = \int sa.\ \int t'.\ f\ (sa\ \#\#\ t')\ \partial T\ (\pi\text{-}Suc\ p\ sa)\ (K\ sa)\ \partial K0\ (p\ [])\ x$
 $\langle proof \rangle$

lemma *nn-integral-T*:

assumes f : $f \in borel\text{-}measurable\ S$
shows $(\int^+ t.\ f\ t\ \partial T\ p\ x) = (\int^+ sa.\ \int^+ t'.\ f\ (sa\ \#\#\ t')\ \partial T\ (\pi\text{-}Suc\ p\ sa)\ (K\ sa)\ \partial K0\ (p\ [])\ x)$
 $\langle proof \rangle$

5.2.3 A Denotational View on the Stochastic Process

Many definitions on MDPs do not rely on the individual traces but only on the distribution of states and actions at each epoch. We define this view on the trace space as the repeated iteration of K_0 and K . It coincides with the definition of T .

primrec $Pn :: ('s,\ 'a)\ pol \Rightarrow 's\ pmf \Rightarrow nat \Rightarrow ('s \times 'a)\ pmf$ **where**
 $Pn\ p\ S0\ 0 = K0\ (p\ [])\ S0$
 $| Pn\ p\ S0\ (Suc\ n) = K0\ (p\ [])\ S0 \gg (\lambda sa.\ Pn\ (\pi\text{-}Suc\ p\ sa)\ (K\ sa)\ n)$
declare $Pn.\text{simps}(2)[simp\ del]$

lemma *Pn-eq-T*:

shows $measure\text{-}pmf\ (Pn\ p\ S0\ n) = distr\ (T\ p\ S0)\ (count\text{-}space\ UNIV)\ (\lambda t.\ t\ !!\ n)$
 $\langle proof \rangle$

The definition of Pn also allows us to easily prove that only enabled actions can occur in the traces of the MDP.

lemma *Pn-in-A*: $is\text{-}policy\ p \Longrightarrow (s,\ a) \in Pn\ p\ S0\ n \Longrightarrow a \in A\ s$

$\langle proof \rangle$

lemma *T-in-A*:

assumes *is-policy* p

shows $AE\ t\ in\ T\ p\ S0.\ snd\ (t\ !!\ n) \in A\ (fst\ (t\ !!\ n))$

$\langle proof \rangle$

5.2.4 State Process

Alongside Pn , we also define the state and action distributions as projections.

definition $Xn\ p\ S0\ n = map\ pmf\ fst\ (Pn\ p\ S0\ n)$

lemma $X0\ [simp]: Xn\ p\ S0\ 0 = S0$

$\langle proof \rangle$

lemma $Xn\text{-}Suc: Xn\ p\ S0\ (Suc\ n) = Pn\ p\ S0\ n \gg K$

$\langle proof \rangle$

lemma $Pn\text{-}markovian\text{-}eq\text{-}Xn\text{-}bind: Pn\ (mk\text{-}markovian\ p)\ S0\ n = K0\ (p\ n)\ (Xn\ (mk\text{-}markovian\ p)\ S0\ n)$

$\langle proof \rangle$

lemma $Xn\text{-}Suc': Xn\ p\ S0\ (Suc\ n) = K0\ (p\ [])\ S0 \gg (\lambda sa.\ Xn\ (\pi\text{-}Suc\ p\ sa)\ (K\ sa)\ n)$

$\langle proof \rangle$

lemma $set\text{-}pmf\text{-}X0\ [simp]: set\text{-}pmf\ (Xn\ p\ S0\ 0) = S0$

$\langle proof \rangle$

lemma $set\text{-}pmf\text{-}PSuc: set\text{-}pmf\ (Pn\ (mk\text{-}markovian\ p)\ S0\ n) =$

$\{(s, a).\ s \in set\text{-}pmf\ (Xn\ (mk\text{-}markovian\ p)\ S0\ n) \wedge a \in p\ n\ s\}$

$\langle proof \rangle$

5.2.5 The Conditional Distribution of Actions

Actions are selected wrt. the whole history of state-action pairs encountered so far. The following definition defines the expected action selection when only the current state is given.

definition $Y\text{-}cond\text{-}X\ p\ S0\ n\ x = map\text{-}pmf\ snd\ (cond\text{-}pmf\ (Pn\ p\ S0\ n)\ \{(s, a).\ s = x\})$

lemma $prob\text{-}K0\text{-}X\ [simp]: measure\text{-}pmf.\ prob\ (K0\ p\ S0)\ \{(s, a).\ s = x\} = pmf\ S0\ x$

$\langle proof \rangle$

lemma $prob\text{-}Pn\text{-}X\ [simp]: measure\text{-}pmf.\ prob\ (Pn\ p\ S0\ n)\ \{(s, a).\ s = x\} = pmf\ (Xn\ p\ S0\ n)\ x$

$\langle proof \rangle$

lemma *pmf-Pn-pair*:

assumes $sa \in set\text{-}pmf (Pn\ p\ S0\ n)$
shows $pmf (Pn\ p\ S0\ n)\ sa = pmf (Y\text{-}cond\text{-}X\ p\ S0\ n\ (fst\ sa))\ (snd\ sa) * pmf (Xn\ p\ S0\ n)\ (fst\ sa)$
 $\langle proof \rangle$

lemma *pmf-Pn*:

assumes $x \in set\text{-}pmf (Xn\ p\ S0\ n)$
shows $pmf (Pn\ p\ S0\ n)\ (x,a) = pmf (Y\text{-}cond\text{-}X\ p\ S0\ n\ x)\ a * pmf (Xn\ p\ S0\ n)\ x$
 $\langle proof \rangle$

lemma *pmf-Y-cond-X*:

assumes $x \in set\text{-}pmf (Xn\ p\ S0\ n)$
shows $pmf (Y\text{-}cond\text{-}X\ p\ S0\ n\ x)\ a = pmf (Pn\ p\ S0\ n)\ (x,a) / pmf (Xn\ p\ S0\ n)\ x$
 $\langle proof \rangle$

lemma *Y-cond-X-0[simp]*:

assumes $x \in set\text{-}pmf\ S0$
shows $Y\text{-}cond\text{-}X\ p\ S0\ 0\ x = p \sqcap x$
 $\langle proof \rangle$

lemma *Y-cond-X-markovian[simp]*:

assumes $h: x \in Xn\ (mk\text{-}markovian\ p)\ S0\ n$
shows $Y\text{-}cond\text{-}X\ (mk\text{-}markovian\ p)\ S0\ n\ x = p\ n\ x$
 $\langle proof \rangle$

lemma *Pn-eq-Xn-Y-cond*: $Pn\ p\ S0\ n = Xn\ p\ S0\ n \ggg (\lambda x. map\text{-}pmf\ (\lambda a. (x, a))\ (Y\text{-}cond\text{-}X\ p\ S0\ n\ x))$
 $\langle proof \rangle$

lemma *Pn-eq-Xn-Y-cond'*:

$Pn\ p\ S0\ n = Xn\ p\ S0\ n \ggg (\lambda s. Y\text{-}cond\text{-}X\ p\ S0\ n\ s \ggg (\lambda a. return\text{-}pmf\ (s,a)))$
 $\langle proof \rangle$

lemma *Pn-markovian-Suc*: $Pn\ (mk\text{-}markovian\ p)\ S0\ (Suc\ n) = Pn\ (mk\text{-}markovian\ p)\ S0\ n \ggg (\lambda sa. K0\ (p\ (Suc\ n))\ (K\ sa))$
 $\langle proof \rangle$

5.2.6 Action Process

The distribution of actions.

definition $Yn\ p\ S0\ n = map\text{-}pmf\ snd\ (Pn\ p\ S0\ n)$

lemma Y_0 : $Y_n p \text{ } S_0 \text{ } 0 = S_0 \gg p \text{ } \square$
 ⟨proof⟩

For markovian policies, the decision rules at each epoch are independent of each other, hence we may express Y_n solely in terms of X_n and the current decision rule.

lemma Y_n -markovian: $Y_n \text{ (mk-markovian } p) \text{ } S_0 \text{ } n = X_n \text{ (mk-markovian } p) \text{ } S_0 \text{ } n \gg p \text{ } n$
 ⟨proof⟩

5.3 Restriction to Markovian Policies

abbreviation as -markovian $p \text{ } S_0 \text{ } n \text{ } x \equiv$
 if $x \in (X_n p \text{ } S_0 \text{ } n)$ then Y -cond- $X \text{ } p \text{ } S_0 \text{ } n \text{ } x$ else return-pmf (SOME
 a. $a \in A \text{ } x$)

For states which cannot occur we choose an arbitrary enabled action, as in this case we cannot make any statements about Y -cond- X (a distribution conditioned on an event with probability 0).

lemma is - Π_{MR} - as -markovian:
assumes p : is -policy p
shows as -markovian $p \text{ } S_0 \in \Pi_{MR}$
 ⟨proof⟩

lemma is -policy- as -markovian: is -policy $p \implies is$ -policy (mk-markovian
 (as -markovian $p \text{ } S_0$))
 ⟨proof⟩

theorem P_n - as -markovian- eq : $P_n \text{ (mk-markovian (} as$ -markovian $p \text{ } S_0)) \text{ } S_0 = P_n p \text{ } S_0$
 ⟨proof⟩

5.4 MDPs without Initial Distribution

From now on, we assume a known, deterministic initial state. All results from the previous discussion carry over as we are now in the special case where the initial state is of the form return-pmf s .

definition $\mathcal{T} \text{ } p \text{ } s \equiv T \text{ } p \text{ (return-pmf } s)$

lemma \mathcal{T} - eq -return-distr: $\mathcal{T} \text{ } p \text{ } s =$
 measure-pmf ($p \text{ } \square \text{ } s$) $\gg (\lambda a. \text{ distr } (T \text{ } (\pi$ -Suc $p \text{ } (s, a)) \text{ } (K \text{ } (s, a)))) \text{ } S$
 ($(\#\#) \text{ } (s, a)$)
 ⟨proof⟩

lemma \mathcal{T} - eq -return:

shows $\mathcal{T} p s = do \{$
 $y \leftarrow measure\text{-}pmf (p \square s);$
 $\omega \leftarrow T (\pi\text{-}Suc p (s,y)) (K (s,y));$
 $return S ((s,y) \#\# \omega)$
 $\}$
 $\langle proof \rangle$

lemma \mathcal{T} -return:

shows $T p S0 = measure\text{-}pmf S0 \gg= T p$
 $\langle proof \rangle$

lemma \mathcal{T} -return-eq:

shows
 $\mathcal{T} p s = do \{$
 $a \leftarrow measure\text{-}pmf (p \square s);$
 $s' \leftarrow measure\text{-}pmf (K (s,a));$
 $w \leftarrow T (\pi\text{-}Suc p (s,a)) (return\text{-}pmf s');$
 $return S ((s,a) \#\# w)$
 $\}$
 $\langle proof \rangle$

lemma \mathcal{T} -eq:

shows $\mathcal{T} p s = do \{$
 $a \leftarrow measure\text{-}pmf (p \square s);$
 $s' \leftarrow measure\text{-}pmf (K (s,a));$
 $w \leftarrow T (\pi\text{-}Suc p (s,a)) s';$
 $return S ((s,a) \#\# w)$
 $\}$
 $\langle proof \rangle$

lemma \mathcal{T} -prob-space[*intro*]: *prob-space* ($\mathcal{T} p s$)

$\langle proof \rangle$

lemma \mathcal{T} -sets[*measurable-cong*]:

$sets (\mathcal{T} p s) = sets S$
 $\langle proof \rangle$

lemma *measurable-ident-Suc'*[*measurable*]:

$(\lambda x. x) \in T (\pi\text{-}Suc p sa) s' \rightarrow_M S$
 $\langle proof \rangle$

lemma *nn-integral-T*:

fixes $f :: ('s \times 'a) \text{ stream} \Rightarrow \text{real}$
assumes $f[\text{measurable}] : f \in \text{borel-measurable } S$
shows $(\int^{+t}. f t \partial \mathcal{T} p s)$
 $= \int^{+a}. \int^{+s'}. \int^{+t'}. f ((s,a) \#\# t') \partial \mathcal{T} (\pi\text{-}Suc p (s,a)) s' \partial K (s,a)$
 $\partial p \square s$
 $\langle proof \rangle$

lemma *integral-T*:

fixes $f :: ('s \times 'a) \text{ stream} \Rightarrow \text{real}$

assumes *f-bounded*: $\bigwedge x. |f x| \leq B$

assumes *f[measurable]*: $f \in \text{borel-measurable } S$

shows $(\int t. f t \partial \mathcal{T} p s)$

$$= \int a. \int s'. \int t'. f ((s,a)\#\#t') \partial \mathcal{T} (\pi\text{-Suc } p (s,a)) s' \partial K (s,a) \partial p$$

\square s

$\langle \text{proof} \rangle$

lemma *integrable-T-bounded[intro]*:

fixes $f :: ('s \times 'a) \text{ stream} \Rightarrow 'd :: \{\text{second-countable-topology, banach}\}$

assumes *f[measurable]*: $f \in \text{borel-measurable } S$

assumes *b: bounded (range f)*

shows *integrable (T p s) f*

$\langle \text{proof} \rangle$

definition $Pn' p s = Pn p$ (*return-pmf s*)

definition $Xn' p s = Xn p$ (*return-pmf s*)

definition $Yn' p s = Yn p$ (*return-pmf s*)

definition $K0' d s \equiv \text{map-pmf } (\lambda a. (s, a)) (d s)$

definition $K\text{-st } d s \equiv d s \gg (\lambda a. K (s, a))$

lemma *pmf-K-st*: $\text{pmf } (K\text{-st } d s) t = \int a. \text{pmf } (K(s, a)) t \partial d s$

$\langle \text{proof} \rangle$

K-st defines the distribution over the successor states for a given decision rule and state. It is mostly useful for markovian policies, as the information which action was selected is lost.

lemma *P0'[simp]*: $Pn' p s 0 = K0' (p []) s$

$\langle \text{proof} \rangle$

lemma *X0'[simp]*: $Xn' p s 0 = \text{return-pmf } s$

$\langle \text{proof} \rangle$

lemma *Pn-return-pmf*: $S0 \gg (\lambda s'. Pn p (\text{return-pmf } s') n) = Pn p S0 n$

$\langle \text{proof} \rangle$

lemma *PSuc'*: $Pn' p s (\text{Suc } n) = K0' (p []) s \gg (\lambda sa. K sa \gg (\lambda s'. Pn' (\pi\text{-Suc } p sa) s' n))$

$\langle \text{proof} \rangle$

lemma *PSuc'-markovian*:

$Pn' (\text{mk-markovian } p) s (\text{Suc } n) = K\text{-st } (p 0) s \gg (\lambda s'. Pn' (\text{mk-markovian } (p \circ \text{Suc})) s' n)$

$\langle \text{proof} \rangle$

lemma *Xn'-Suc*: $Xn' p s (\text{Suc } n) = Pn' p s n \gg K$

<proof>

lemma *Xn'-Pn'*: $Xn' p s n = \text{map-pmf fst } (Pn' p s n)$
<proof>

lemma *Suc-Xn'*: $Xn' p s (\text{Suc } n) = p [] s \gg= (\lambda a. K (s,a) \gg= (\lambda s'. Xn' (\pi\text{-Suc } p (s,a)) s' n))$
<proof>

lemma *Suc-Xn'-markovian*:
 $Xn' (\text{mk-markovian } p) s (\text{Suc } n) = K\text{-st } (p 0) s \gg= (\lambda s'. Xn' (\text{mk-markovian } (\lambda n. p (\text{Suc } n))) s' n)$
<proof>

lemma *Xn'-split*: $Xn' (\text{mk-markovian } p) s (n + m) = Xn' (\text{mk-markovian } p) s n \gg= (\lambda s. Xn' (\text{mk-markovian } (\lambda i. p (i + n))) s m)$
<proof>

lemma *Yn'-markovian*: $Yn' (\text{mk-markovian } p) s n = Xn' (\text{mk-markovian } p) s n \gg= p n$
<proof>

lemma *Pn'-markovian-eq-Xn'-bind*: $Pn' (\text{mk-markovian } p) s n = Xn' (\text{mk-markovian } p) s n \gg= K0' (p n)$
<proof>

lemma *Pn'-eq-T*: $\text{measure-pmf } (Pn' p s n) = \text{distr } (\mathcal{T} p s) (\text{count-space UNIV}) (\lambda t. t !! n)$
<proof>

end
end

theory *MDP-reward*
imports
 Bounded-Functions
 MDP-reward-Util
 Blinfun-Util
 MDP-disc
begin

6 Markov Decision Processes with Rewards

locale *MDP-reward = discrete-MDP* *A K*
for
 A **and**
 K :: 's :: countable × 'a :: countable ⇒ 's pmf +
fixes

```

    r :: ('s × 'a) ⇒ real and
    l :: real
assumes
    zero-le-disc [simp]: 0 ≤ l and
    disc-lt-one [simp]: l < 1 and
    r-bounded: bounded (range r)
begin

```

This extension to the basic MDPs is formalized with another locale. It assumes the existence of a reward function r which takes a state-action pair to a real number. We assume that the function is bounded r -bounded.

Furthermore, we fix a discounting factor l , where $0 \leq l \wedge l < 1$.

6.1 Basic Properties

```

lemma r-bfun: r ∈ bfun
  ⟨proof⟩

```

```

lemma r-bounded': bounded (r ' X)
  ⟨proof⟩

```

```

lemma abs-disc-eq[simp]: |l ^ i * x| = l ^ i * |x|
  ⟨proof⟩

```

```

definition r_M = (⊔ sa. |r sa|)

```

```

lemma abs-r-le-r_M: |r sa| ≤ r_M
  ⟨proof⟩

```

```

lemma abs-r_M-eq-r_M [simp]: |r_M| = r_M
  ⟨proof⟩

```

```

lemma r_M-nonneg: 0 ≤ r_M
  ⟨proof⟩

```

6.2 Summability

```

lemma summable-disc [intro, simp]: summable (λi. l ^ i * x)
  ⟨proof⟩

```

```

lemma summable-r-disc[intro, simp]:
  summable (λi. |l ^ i * r (sa i)|)
  summable (λi. l ^ i * |r (sa i)|)
  summable (λi. l ^ i * r (sa i))
  ⟨proof⟩

```

6.3 Reward over a Trace

abbreviation $\nu\text{-trace-fin } t N \equiv \sum_{i < N}. l \hat{\ } i * r (t !! i)$

abbreviation $\nu\text{-trace } t \equiv \sum_{i}. l \hat{\ } i * r (t !! i)$

lemma $\text{abs-}\nu\text{-trace-le}$: $|\nu\text{-trace } t| \leq (\sum_{i}. l \hat{\ } i * r_M)$

$\langle \text{proof} \rangle$

lemma $\text{abs-}\nu\text{-trace-fin-le}$: $|\nu\text{-trace-fin } t N| \leq (\sum_{i < N}. l \hat{\ } i * r_M)$

$\langle \text{proof} \rangle$

lemma $\text{measurable-suminf-reward}$ [*measurable*]: $\nu\text{-trace} \in \text{borel-measurable } S$

$\langle \text{proof} \rangle$

lemma $\text{integrable-}\nu\text{-trace-fin}$: $\text{integrable } (\mathcal{T} \ p \ s) (\lambda t. \nu\text{-trace-fin } t N)$

$\langle \text{proof} \rangle$

lemma $\text{integrable-}\nu\text{-trace}$: $\text{integrable } (\mathcal{T} \ p \ s) \ \nu\text{-trace}$

$\langle \text{proof} \rangle$

6.4 Integrals over Rewards

lemma measurable-r-nth [*measurable*]: $(\lambda t. r (t !! i)) \in \text{borel-measurable } S$

S

$\langle \text{proof} \rangle$

lemma integrable-r-nth [*simp*]: $\text{integrable } (\mathcal{T} \ p \ s) (\lambda t. r (t !! i))$

$\langle \text{proof} \rangle$

lemma $\text{expectation-abs-r-le}$: $\text{measure-pmf.expectation } d (\lambda a. |r (s, a)|)$

$\leq r_M$

$\langle \text{proof} \rangle$

6.5 Expected Total Discounted Reward

context

fixes $p :: ('s, 'a) \text{ pol}$

begin

6.6 Expected Finite-Horizon Discounted Reward

definition $\nu\text{-fin } n \ s = \int t. \nu\text{-trace-fin } t \ n \ \partial \mathcal{T} \ p \ s$

lemma $\text{abs-}\nu\text{-fin-le}$: $|\nu\text{-fin } N \ s| \leq (\sum_{i < N}. l \hat{\ } i * r_M)$

$\langle \text{proof} \rangle$

lemma $\nu\text{-fin-Suc}$ [*simp*]: $\nu\text{-fin } (\text{Suc } n) \ s = \nu\text{-fin } n \ s + l \hat{\ } n * \int t. r$

$(t !! n) \ \partial \mathcal{T} \ p \ s$

$\langle \text{proof} \rangle$

lemma ν -fin-zero[simp]: ν -fin 0 s = 0
⟨proof⟩

lemma ν -fin-eq-Pn: ν -fin n s = ($\sum i < n. l \hat{i} * \text{measure-pmf.expectation}$
(Pn' p s i) r)
⟨proof⟩

definition ν s = lim ($\lambda n. \nu$ -fin n s)

lemma ν -eq-lim: ν s = lim ($\lambda n. \nu$ -fin n s)
⟨proof⟩

lemma ν -eq- ν -trace: ν s = $\int t. \nu$ -trace t $\partial \mathcal{T}$ p s
⟨proof⟩

lemma abs- ν -le: $|\nu$ s| \leq ($\sum i. l \hat{i} * r_M$)
⟨proof⟩

lemma ν -le: ν s \leq ($\sum i. l \hat{i} * r_M$)
⟨proof⟩

lemma ν -eq-Pn: ν s = ($\sum i. l \hat{i} * \text{measure-pmf.expectation}$ (Pn' p s
i) r)
⟨proof⟩

lemma ν -bfun: $\nu \in \text{bfun}$
⟨proof⟩

lemma ν -fin-bfun: ($\lambda s. \nu$ -fin N s) $\in \text{bfun}$
⟨proof⟩

lift-definition $\nu_b :: 's \Rightarrow_b \text{real}$ is ν
⟨proof⟩

lemma norm- ν -le: norm $\nu_b \leq r_M / (1-l)$
⟨proof⟩

end

lemma ν -as-markovian: ν (mk-markovian (as-markovian p (return-pmf
s))) s = ν p s
⟨proof⟩

6.6.1 Optimal Reward

definition ν -MD s $\equiv \bigsqcup p \in \Pi_{MD}. \nu$ (mk-markovian-det p) s

definition ν -opt s $\equiv \bigsqcup p \in \Pi_{HR}. \nu$ p s

lemma ν -le- ν -opt [intro]:
assumes *is-policy* p
shows $\nu p s \leq \nu$ -opt s
 \langle proof \rangle

lemma ν -opt-eq-MR: ν -opt $s = (\bigsqcup p \in \Pi_{MR}. \nu (mk$ -markovian $p) s)$
 \langle proof \rangle

lemma ν -opt-bfun: ν -opt \in bfun
 \langle proof \rangle

lift-definition ν_b -opt :: ' $s \Rightarrow_b$ real is ν -opt
 \langle proof \rangle

6.7 Reward of a Decision Rule

context
fixes $d :: ('s, 'a) dec$
begin
abbreviation r -dec $s \equiv \int a. r (s, a) \partial d s$

lemma *abs-r-dec-le*: $|r$ -dec $s| \leq r_M$
 \langle proof \rangle

lemma *r-dec-eq-r-K0*: r -dec $s = measure$ -pmf.*expectation* ($K0'$ $d s$) r
 \langle proof \rangle

lemma *r-dec-bfun*: r -dec \in bfun
 \langle proof \rangle

lift-definition r -dec _{b} :: ' $s \Rightarrow_b$ real is r -dec
 \langle proof \rangle

lemma *norm-r-dec-le*: $norm$ r -dec _{b} $\leq r_M$
 \langle proof \rangle

end

lemma *r-dec-det [simp]*: r -dec (mk -dec-det d) $s = r (s, d s)$
 \langle proof \rangle

declare r -dec _{b} .*rep-eq[simp]* bfun.*Bfun-inverse[simp]*

6.8 Push-Forward of a Function Through the MDP

lemma *norm-l-pow-eq[simp]*: $norm (l \hat{t} *_R F) = l \hat{t} * norm F$
 \langle proof \rangle

lemma *summable-norm-disc-I [intro]*:
assumes *summable* ($\lambda t. (l \hat{t} * norm F)$)

shows *summable* ($\lambda t. \text{norm } (l \hat{=} t *_{\mathbb{R}} F)$)
<proof>

lemma *summable-norm-disc-I'* [intro]:
assumes *summable* ($\lambda t. (l \hat{=} t * \text{norm } (F t))$)
shows *summable* ($\lambda t. \text{norm } (l \hat{=} t *_{\mathbb{R}} F t)$)
<proof>

lemma *summable-discI* [intro]:
assumes *bounded* (*range* F)
shows *summable* ($\lambda t. l \hat{=} t * \text{norm } (F t)$)
<proof>

lemma *summable-disc-reward* [intro]:
assumes *bounded* (*range* ($F :: \text{nat} \Rightarrow 'b :: \text{banach}$))
shows *summable* ($\lambda t. l \hat{=} t *_{\mathbb{R}} (F t)$)
<proof>

context
fixes $p :: \text{nat} \Rightarrow ('s, 'a) \text{dec}$

begin

definition $\mathcal{P}_X \ n = \text{push-exp } (\lambda s. Xn' (\text{mk-markovian } p) \ s \ n)$

lemma $\mathcal{P}_X\text{-}0$ [simp]: $\mathcal{P}_X \ 0 = \text{id}$
<proof>

lemma $\mathcal{P}_X\text{-bounded-linear}$ [simp]: *bounded-linear* ($\mathcal{P}_X \ t$)
<proof>

lemma $\text{norm-}\mathcal{P}_X$ [simp]: *onorm* ($\mathcal{P}_X \ t$) = 1
<proof>

lemma $\text{norm-}\mathcal{P}_X\text{-apply}$ [simp]: $\text{norm } (\mathcal{P}_X \ n \ x) \leq \text{norm } x$
<proof>

lemma $\mathcal{P}_X\text{-bound-r}$: $\text{norm } (\mathcal{P}_X \ t \ (r\text{-dec}_b \ (p \ t))) \leq r_M$
<proof>

lemma $\mathcal{P}_X\text{-bounded-r}$: *bounded* (*range* ($\lambda t. (\mathcal{P}_X \ t \ (r\text{-dec}_b \ (p \ t))))$)
<proof>

lemma *summable-norm-disc-reward'* [simp]:
shows *summable* ($\lambda t. l \hat{=} t * \text{norm } (\mathcal{P}_X \ t \ (r\text{-dec}_b \ (p \ t)))$)
<proof>

lemma *summable-disc-reward-}\mathcal{P}_X [simp]:
shows *summable* ($\lambda t. l \hat{=} t *_{\mathbb{R}} \mathcal{P}_X \ t \ (r\text{-dec}_b \ (p \ t))$)
*<proof>**

lemma *disc-reward-tendsto*:

$(\lambda n. \sum t < n. l \hat{t} *_R \mathcal{P}_X t (r\text{-dec}_b (p t))) \longrightarrow (\sum t. l \hat{t} *_R \mathcal{P}_X t (r\text{-dec}_b (p t)))$
 $\langle \text{proof} \rangle$

end

lemma $\mathcal{P}_X\text{-Suc}$: $\mathcal{P}_X p (Suc n) v = \text{push-exp } (K\text{-st } (p 0)) ((\mathcal{P}_X (\lambda n. p (Suc n)) n) v)$
 $\langle \text{proof} \rangle$

lemma $\mathcal{P}_X\text{-Suc}'$: $\mathcal{P}_X p (Suc n) v = \mathcal{P}_X p n (\text{push-exp } (K\text{-st } (p n)) v)$
 $\langle \text{proof} \rangle$

lemma $\mathcal{P}_X\text{-sconst}$: $\mathcal{P}_X (\lambda \cdot. p) n = (\text{push-exp } (K\text{-st } p)) \hat{\sim} n$
 $\langle \text{proof} \rangle$

lemma *norm-P-n[simp]*: $\text{onorm } (\text{push-exp } (K\text{-st } p) \hat{\sim} n) = 1$
 $\langle \text{proof} \rangle$

lemma *summable-norm-bfun-disc*: $\text{summable } (\lambda t. l \hat{t} * \text{norm } (\text{apply-bfun } f t))$
 $\langle \text{proof} \rangle$

lemma *summable-bfun-disc [simp]*: $\text{summable } (\lambda t. l \hat{t} * (\text{apply-bfun } f t))$
 $\langle \text{proof} \rangle$

lemma *summable-disc-norm*: $\text{summable } (\lambda x. l \hat{x} * \text{norm } c)$
 $\langle \text{proof} \rangle$

lemma *norm-bfun-disc-le*: $\text{norm } f \leq B \implies (\sum x. l \hat{x} * \text{norm } (\text{apply-bfun } f x)) \leq (\sum x. l \hat{x} * B)$
 $\langle \text{proof} \rangle$

lemma *norm-bfun-disc-le'*: $\text{norm } f \leq B \implies (\sum x. l \hat{x} * (\text{apply-bfun } f x)) \leq (\sum x. l \hat{x} * B)$
 $\langle \text{proof} \rangle$

lemma *sum-disc-lim-l*: $(\sum x. l \hat{x} * B) = B / (1-l)$
 $\langle \text{proof} \rangle$

lemma *sum-disc-bound*: $(\sum x. l \hat{x} * \text{apply-bfun } f x) \leq (\text{norm } f) / (1-l)$
 $\langle \text{proof} \rangle$

lemma *sum-disc-bound'*:
fixes $f :: \text{nat} \Rightarrow 'b \Rightarrow_b \text{real}$

assumes $h: \forall n. \text{norm } (f n) \leq B$
shows $\text{norm } (\sum x. l\hat{x} *_R f x) \leq B / (1-l)$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-le-opt}$ [intro]: $p \in \Pi_{HR} \implies \nu_b p \leq \nu_b\text{-opt}$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-le-opt-MD}$ [intro]: $p \in \Pi_{MD} \implies \nu_b (\text{mk-markovian-det } p) \leq \nu_b\text{-opt}$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-le-opt-DD}$ [intro]: $\text{is-dec-det } d \implies \nu_b (\text{mk-stationary-det } d) \leq \nu_b\text{-opt}$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-le-opt-DR}$ [intro]: $\text{is-dec } d \implies \nu_b (\text{mk-stationary } d) \leq \nu_b\text{-opt}$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-opt-eq-MR}$: $\nu_b\text{-opt } s = (\bigsqcup p \in \Pi_{MR}. \nu_b (\text{mk-markovian } p) s)$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-as-markovian}$: $\nu_b (\text{mk-markovian } (\text{as-markovian } p (\text{return-pmf } s))) s = \nu_b p s$
 $\langle \text{proof} \rangle$

lemma $\nu\text{-elem}$: $\nu (\text{mk-markovian } p) s = (\sum i. l\hat{i} * \mathcal{P}_X p i (\text{r-dec}_b (p i)) s)$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-eq-}\mathcal{P}_X$: $\nu_b (\text{mk-markovian } p) = (\sum i. l\hat{i} *_R \mathcal{P}_X p i (\text{r-dec}_b (p i)))$
 $\langle \text{proof} \rangle$

lemma $\nu\text{-eq-}\mathcal{P}_X$: $\nu (\text{mk-markovian } p) = (\sum i. l\hat{i} *_R \mathcal{P}_X p i (\text{r-dec}_b (p i)))$
 $\langle \text{proof} \rangle$

$\mathcal{P}_1 d v$ defines for each state the expected value of v after taking a single step in the MDP according to the decision rule d .

context

fixes $d :: ('s, 'a) \text{dec}$

begin

lift-definition $\mathcal{P}_1 :: ('s \Rightarrow_b \text{real}) \Rightarrow_L ('s \Rightarrow_b \text{real})$ **is** $\text{push-exp } (K\text{-st } d)$

$\langle \text{proof} \rangle$

lemma $\mathcal{P}_1\text{-pow}$: $\text{blifun-apply } (\mathcal{P}_1 \hat{\sim} n) = \text{blifun-apply } \mathcal{P}_1 \hat{\sim} n$

$\langle proof \rangle$

lemma \mathcal{P}_X -const: $\mathcal{P}_X (\lambda \cdot d) n = \mathcal{P}_1 \hat{\sim} n$
 $\langle proof \rangle$

lemma norm- \mathcal{P}_1 [simp]: norm $\mathcal{P}_1 = 1$
 $\langle proof \rangle$

lemma norm- \mathcal{P}_1 -pow [simp]: norm $(\mathcal{P}_1 \hat{\sim} t) = 1$
 $\langle proof \rangle$

lemma norm- \mathcal{P}_1 -l-less: norm $(l *_R \mathcal{P}_1) < 1$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -pos: $0 \leq u \implies 0 \leq \mathcal{P}_1 u$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -n-pos: $0 \leq u \implies 0 \leq (\mathcal{P}_1 \hat{\sim} n) u$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -n-disc-pos: $0 \leq u \implies 0 \leq (l \hat{\sim} n *_R \mathcal{P}_1 \hat{\sim} n) u$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -sum-pos: $0 \leq u \implies 0 \leq (\sum_{t \leq n}. l \hat{\sim} t *_R (\mathcal{P}_1 \hat{\sim} t)) u$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -sum-ge:
assumes $0 \leq u$
shows $u \leq (\sum_{t \leq n}. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t) u$
 $\langle proof \rangle$

lemma disc- \mathcal{P}_1 -tendsto: $(\lambda n. (\sum_{t \leq n}. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t)) \longrightarrow (\sum t. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t)$
 $\langle proof \rangle$

lemma disc- \mathcal{P}_1 -lim: $\lim (\lambda n. (\sum_{t \leq n}. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t)) = (\sum t. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t)$
 $\langle proof \rangle$

lemma convergent-disc- \mathcal{P}_1 : convergent $(\lambda n. (\sum_{t \leq n}. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t))$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -suminf-ge:
assumes $0 \leq u$ shows $u \leq (\sum t. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t) u$
 $\langle proof \rangle$

lemma \mathcal{P}_1 -suminf-pos:
assumes $0 \leq u$
shows $0 \leq (\sum t. l \hat{\sim} t *_R \mathcal{P}_1 \hat{\sim} t) u$

<proof>

lemma *lemma-6-1-2-b:*

assumes $v \leq u$

shows $(\sum t. l \hat{t} *_R \mathcal{P}_1 \hat{\sim} t) v \leq (\sum t. l \hat{t} *_R \mathcal{P}_1 \hat{\sim} t) u$

<proof>

6.9 The Bellman Operator L

definition $L v \equiv r\text{-dec}_b d + l *_R \mathcal{P}_1 v$

lemma *norm-L-le:* $\text{norm } (L v) \leq r_M + l * \text{norm } v$

<proof>

lemma *abs-L-le:* $|L v s| \leq r_M + l * \text{norm } v$

<proof>

lemma *ν -stationary:* $\nu_b (\text{mk-stationary } d) = (\sum t. l \hat{t} *_R (\mathcal{P}_1 \hat{\sim} t))$

$(r\text{-dec}_b d)$

<proof>

end

lemma *\mathcal{P}_1 -eq- \mathcal{P}_X -one:* $\text{blinfun-apply } (\mathcal{P}_1 (p \ 0)) = \mathcal{P}_X p \ 1$

<proof>

The value of a markovian policy can be expressed in terms of L .

lemma *ν -step:* $\nu_b (\text{mk-markovian } p) = L (p \ 0) (\nu_b (\text{mk-markovian } (\lambda n. p (\text{Suc } n))))$

<proof>

lemma *L - ν -fix:* $\nu_b (\text{mk-stationary } d) = L d (\nu_b (\text{mk-stationary } d))$

<proof>

lemma *L -fix- ν :*

assumes $L p v = v$

shows $v = \nu_b (\text{mk-stationary } p)$

<proof>

lemma *L - ν -fix-iff:* $L d v = v \longleftrightarrow v = \nu_b (\text{mk-stationary } d)$

<proof>

lemma *apply-bfun-bounded-above* [*simp, intro*]:

fixes $f :: 'c \Rightarrow_b \text{real}$

shows $\text{bounded } (f \text{ ' } X)$

<proof>

lemma *apply-bfun-bdd-above*[*simp*, *intro*]:

fixes $f :: 'c \Rightarrow_b \text{real}$

shows *bdd-above* ($f \text{ ' } X$)

<proof>

lemma *L-bounded*[*simp*, *intro*]: *bounded* (*range* ($\lambda p. L p v s$))

<proof>

lemma *L-bounded'*[*simp*, *intro*]: *bounded* ($(\lambda p. L p v s) \text{ ' } X$)

<proof>

lemma *L-bdd-above*[*simp*, *intro*]: *bdd-above* ($(\lambda p. L p v s) \text{ ' } X$)

<proof>

6.10 Optimality Equations

definition $\mathcal{L} (v :: 's \Rightarrow_b \text{real}) s = (\bigsqcup d \in D_R. L d v s)$

lemma *L-bfun*: $\mathcal{L} v \in \text{bfun}$

<proof>

lift-definition $\mathcal{L}_b :: ('s \Rightarrow_b \text{real}) \Rightarrow 's \Rightarrow_b \text{real}$ **is** \mathcal{L}

<proof>

lemma *L-le-L_b*: *is-dec* $d \implies L d v \leq \mathcal{L}_b v$

<proof>

6.11 Monotonicity

lemma *P₁-mono*[*intro*]: $a \leq b \implies \mathcal{P}_1 p a \leq \mathcal{P}_1 p b$

<proof>

lemma *P_X-mono*[*intro*]: $a \leq b \implies \mathcal{P}_X p n a \leq \mathcal{P}_X p n b$

<proof>

lemma *L-mono*: $u \leq v \implies L d u \leq L d v$

<proof>

lemma *L_b-mono*:

assumes $u \leq v$ **shows** $\mathcal{L}_b u \leq \mathcal{L}_b v$

<proof>

lemma *step-mono*:

assumes $\mathcal{L}_b v \leq v$ $d \in D_R$

shows $L d v \leq v$

<proof>

lemma *step-mono-elem*:

assumes $v \leq \mathcal{L}_b v e > 0$
shows $\exists d \in D_R. v \leq L d v + e *_R 1$
 $\langle \text{proof} \rangle$

lemma $p\text{-}n\text{-}\pi\text{-}MD[\text{intro}]$: $p \in \Pi_{MD} \implies p n \in D_D$
 $\langle \text{proof} \rangle$

lemma $p\text{-}n\text{-}\pi\text{-}MR[\text{intro}]$: $p \in \Pi_{MR} \implies p n \in D_R$
 $\langle \text{proof} \rangle$

lemma $\mathcal{P}_X\text{-}Suc\text{-}n\text{-}elem$: $\mathcal{P}_X p n (\mathcal{P}_1 (p n) v) = \mathcal{P}_X p (Suc n) v$
 $\langle \text{proof} \rangle$

lift-definition $\nu_b\text{-}fin$:: $('s, 'a) \text{pol} \Rightarrow \text{nat} \Rightarrow 's \Rightarrow_b \text{real}$ **is** $\nu\text{-}fin$
 $\langle \text{proof} \rangle$

lemma $\nu\text{-}fin\text{-}elem$: $\nu\text{-}fin (mk\text{-}markovian p) n s = (\sum i < n. l\hat{i} * \mathcal{P}_X p i (r\text{-}dec_b (p i)) s)$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-}fin\text{-}eq\text{-}\mathcal{P}_X$: $\nu_b\text{-}fin (mk\text{-}markovian p) n = (\sum i < n. l\hat{i} *_R \mathcal{P}_X p i (r\text{-}dec_b (p i)))$
 $\langle \text{proof} \rangle$

lemma $\nu\text{-}fin\text{-}eq\text{-}\mathcal{P}_X$: $\nu\text{-}fin (mk\text{-}markovian p) n = (\sum i < n. l\hat{i} *_R \mathcal{P}_X p i (r\text{-}dec_b (p i)))$
 $\langle \text{proof} \rangle$

abbreviation $\mathcal{P}_d p n v \equiv l\hat{n} *_R \mathcal{P}_X p n v$

lemma $\mathcal{P}_X\text{-}L\text{-}le$:
assumes $\mathcal{L}_b v \leq v p \in \Pi_{MR}$
shows $\mathcal{P}_X p n (L (p n) v) \leq \mathcal{P}_X p n v$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-}fin\text{-}tendsto\text{-}\nu_b$: $(\nu_b\text{-}fin (mk\text{-}markovian p)) \longrightarrow \nu_b (mk\text{-}markovian p)$
 $\langle \text{proof} \rangle$

lemma $\mathcal{P}_d\text{-}lim$: $(\lambda n. (\mathcal{P}_d p n v)) \longrightarrow 0$
 $\langle \text{proof} \rangle$

lemma $\mathcal{P}_1\text{-}bfun\text{-}one$ [simp]: $\mathcal{P}_1 p 1 = 1$
 $\langle \text{proof} \rangle$

lemma $\mathcal{P}_1\text{-}pow\text{-}bfun\text{-}one$ [simp]: $((\mathcal{P}_1 p) \hat{\sim} t) 1 = 1$
 $\langle \text{proof} \rangle$

6.12 Properties of Solutions of the Optimality Equations

lemma *\mathcal{L} -dec-ge-opt*:
 assumes $\mathcal{L}_b v \leq v$
 shows $\nu_b\text{-opt} \leq v$
 $\langle\text{proof}\rangle$

lemma *\mathcal{L} -inc-le-opt*:
 assumes $v \leq \mathcal{L}_b v$
 shows $v \leq \nu_b\text{-opt}$
 $\langle\text{proof}\rangle$

lemma *\mathcal{L} -fix-imp-opt*:
 assumes $v = \mathcal{L}_b v$
 shows $v = \nu_b\text{-opt}$
 $\langle\text{proof}\rangle$

lemma *bounded-P*: *bounded* (\mathcal{P}_1 ‘ X)
 $\langle\text{proof}\rangle$

6.13 Solutions to the Optimality Equation

6.13.1 \mathcal{L}_b and L are Contraction Mappings

declare *bounded-apply-blinfun*[*intro*] *bounded-apply-bfun*'[*intro*]

lemma *contraction- \mathcal{L}* : *dist* ($\mathcal{L}_b v$) ($\mathcal{L}_b u$) $\leq l * \text{dist } v u$
 $\langle\text{proof}\rangle$

lemma *is-contraction- \mathcal{L}* : *is-contraction* \mathcal{L}_b
 $\langle\text{proof}\rangle$

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

lemma *is-contraction-L*: *is-contraction* ($L p$)
 $\langle\text{proof}\rangle$

6.13.2 Existence of a Fixpoint of \mathcal{L}_b

lemma *\mathcal{L}_b -conv*:
 $\exists! v. \mathcal{L}_b v = v$
 $(\lambda n. (\mathcal{L}_b \widehat{\sim} n) v) \longrightarrow (\text{THE } v. \mathcal{L}_b v = v)$
 $\langle\text{proof}\rangle$

lemma *\mathcal{L}_b -fix-iff-opt* [*simp*]: $\mathcal{L}_b v = v \longleftrightarrow v = \nu_b\text{-opt}$
 $\langle\text{proof}\rangle$

lemma ν_b -opt-fix: ν_b -opt = (THE v. $\mathcal{L}_b v = v$)
 ⟨proof⟩

lemma \mathcal{L}_b -opt [simp]: $\mathcal{L}_b \nu_b$ -opt = ν_b -opt
 ⟨proof⟩

lemma \mathcal{L}_b -lim: $(\lambda n. (\mathcal{L}_b \overset{\sim}{\sim} n) v) \longrightarrow \nu_b$ -opt
 ⟨proof⟩

lemma thm-6-2-6: $\nu_b p = \nu_b$ -opt $\longleftrightarrow \mathcal{L}_b (\nu_b p) = \nu_b p$
 ⟨proof⟩

lemma thm-6-2-6': $\nu p = \nu$ -opt $\longleftrightarrow \mathcal{L}_b (\nu_b p) = \nu_b p$
 ⟨proof⟩

6.14 Existence of Optimal Policies

definition ν -improving v d $\longleftrightarrow (\forall s. \text{is-arg-max } (\lambda d. (L d v) s) (\lambda d. d \in D_R) d)$

lemma ν -improving-iff: ν -improving v d $\longleftrightarrow d \in D_R \wedge (\forall d' \in D_R. \forall s. L d' v s \leq L d v s)$
 ⟨proof⟩

lemma ν -improving-D-MR[dest]: ν -improving v d $\implies d \in D_R$
 ⟨proof⟩

lemma ν -improving-ge: ν -improving v d $\implies d' \in D_R \implies L d' v s \leq L d v s$
 ⟨proof⟩

lemma ν -improving-imp- \mathcal{L}_b : ν -improving v d $\implies \mathcal{L}_b v = L d v$
 ⟨proof⟩

lemma \mathcal{L}_b -imp- ν -improving:
 assumes $d \in D_R \mathcal{L}_b v = L d v$
 shows ν -improving v d
 ⟨proof⟩

lemma ν -improving-alt:
 assumes $d \in D_R$
 shows ν -improving v d $\longleftrightarrow \mathcal{L}_b v = L d v$
 ⟨proof⟩

definition ν -conserving d = ν -improving (ν_b -opt) d

lemma ν -conserving-iff: ν -conserving d $\longleftrightarrow d \in D_R \wedge (\forall d' \in D_R. \forall s. L d' \nu_b$ -opt s $\leq L d \nu_b$ -opt s)

$\langle proof \rangle$

lemma ν -conserving-ge: ν -conserving $d \implies d' \in D_R \implies L d' \nu_b\text{-opt}$
 $s \leq L d \nu_b\text{-opt} s$
 $\langle proof \rangle$

lemma ν -conserving-imp- \mathcal{L}_b [simp]: ν -conserving $d \implies L d \nu_b\text{-opt} =$
 $\nu_b\text{-opt}$
 $\langle proof \rangle$

lemma \mathcal{L}_b -imp- ν -conserving:
assumes $d \in D_R$ $\mathcal{L}_b \nu_b\text{-opt} = L d \nu_b\text{-opt}$
shows ν -conserving d
 $\langle proof \rangle$

lemma ν -conserving-alt:
assumes $d \in D_R$
shows ν -conserving $d \iff \mathcal{L}_b \nu_b\text{-opt} = L d \nu_b\text{-opt}$
 $\langle proof \rangle$

lemma ν -conserving-alt':
assumes $d \in D_R$
shows ν -conserving $d \iff L d \nu_b\text{-opt} = \nu_b\text{-opt}$
 $\langle proof \rangle$

6.14.1 Conserving Decision Rules are Optimal

theorem ex -improving-imp-conserving:
assumes $\bigwedge v. \exists d. \nu$ -improving v ($mk\text{-dec-det } d$)
shows $\exists d. \nu$ -conserving ($mk\text{-dec-det } d$)
 $\langle proof \rangle$

theorem $conserving$ -imp-opt[simp]:
assumes ν -conserving ($mk\text{-dec-det } d$)
shows ν_b ($mk\text{-stationary-det } d$) = $\nu_b\text{-opt}$
 $\langle proof \rangle$

lemma $conserving$ -imp-opt':
assumes $\exists d. \nu$ -conserving ($mk\text{-dec-det } d$)
shows $\exists d \in D_D. (\nu_b$ ($mk\text{-stationary-det } d$)) = $\nu_b\text{-opt}$
 $\langle proof \rangle$

theorem $improving$ -att-imp-det-opt:
assumes $\bigwedge v. \exists d. \nu$ -improving v ($mk\text{-dec-det } d$)
shows $\nu_b\text{-opt } s = (\bigsqcup d \in D_D. \nu_b$ ($mk\text{-stationary-det } d$)) s
 $\langle proof \rangle$

6.14.2 Bellman Operator for Single Actions

abbreviation $L_a a v s \equiv r (s, a) + l * \text{measure-pmf.expectation} (K (s,a)) v$

lemma $L_a\text{-le}$:

fixes $v :: 's \Rightarrow_b \text{real}$

shows $|L_a a v s| \leq r_M + l * \text{norm } v$

$\langle \text{proof} \rangle$

lemma $L_a\text{-bounded}$:

$\text{bounded} (\text{range} (\lambda a. L_a a (\text{apply-bfun } v) s))$

$\langle \text{proof} \rangle$

lemma $L_a\text{-int}$:

fixes $d :: 'a \text{ pmf}$ **and** $v :: 's \Rightarrow_b \text{real}$

shows $(\int a. L_a a v s \partial d) = (\int a. r (s, a) \partial d) + l * \int a. \int s'. v s' \partial K (s, a) \partial d$

$\langle \text{proof} \rangle$

lemma $L\text{-eq-}L_a$: $L d v s = \text{measure-pmf.expectation} (d s) (\lambda a. L_a a v s)$

$\langle \text{proof} \rangle$

lemma $L\text{-eq-}L_a\text{-det}$: $L (\text{mk-dec-det } d) v s = L_a (d s) v s$

$\langle \text{proof} \rangle$

6.14.3 Equivalences involving \mathcal{L}_b

lemma $SUP\text{-step-MR-eq}$:

$(\bigsqcup d \in D_R. L d v s) = (\bigsqcup pa \in \{pa. \text{set-pmf } pa \subseteq A s\}. (\int a. L_a a v s \partial \text{measure-pmf } pa))$

$\langle \text{proof} \rangle$

lemma $\mathcal{L}_b\text{-sup-att-det}$:

assumes $d \in D_R \mathcal{L}_b v = L d v$

shows $\exists d' \in D_D. \mathcal{L}_b v = L (\text{mk-dec-det } d') v$

$\langle \text{proof} \rangle$

lemma $SUP\text{-step-det-eq}$: $(\bigsqcup d \in D_D. L (\text{mk-dec-det } d) v s) = (\bigsqcup a \in A s. L_a a v s)$

$\langle \text{proof} \rangle$

lemma $\text{integrable-}L_a$: $\text{integrable} (\text{measure-pmf } x) (\lambda a. L_a a (\text{apply-bfun } v) s)$

$\langle \text{proof} \rangle$

lemma $SUP\text{-}L_a\text{-eq-det}$:

fixes $v :: 's \Rightarrow_b \text{real}$

shows $(\bigsqcup p \in \{p. \text{set-pmf } p \subseteq A s\}. \int a. L_a a v s \partial \text{measure-pmf } p) =$

$(\bigsqcup_{a \in A} s. L_a a v s)$
 $\langle \text{proof} \rangle$

lemma \mathcal{L} -eq-SUP-det: $\mathcal{L} v s = (\bigsqcup_{d \in D_D}. L (mk\text{-dec-det } d) v s)$
 $\langle \text{proof} \rangle$

lemma \mathcal{L}_b -eq-SUP-det: $\mathcal{L}_b v s = (\bigsqcup_{d \in D_D}. L (mk\text{-dec-det } d) v s)$
 $\langle \text{proof} \rangle$

6.14.4 Deterministic Decision Rules are Optimal

lemma *opt-imp-opt-dec-det*:
assumes $p \in \Pi_{HR} \nu_b p = \nu_b\text{-opt}$
shows $\exists d \in D_D. \nu_b (mk\text{-stationary } (mk\text{-dec-det } d)) = \nu_b\text{-opt}$
 $\langle \text{proof} \rangle$

6.14.5 Optimal Decision Rules for Finite Action Spaces

lemma *thm-6-2-10*:
assumes $\bigwedge s. \text{finite } (A s)$
shows $\exists d \in D_D. \nu_b\text{-opt} = \nu_b (mk\text{-stationary-det } d)$
 $\langle \text{proof} \rangle$

6.14.6 Existence of Epsilon-Optimal Policies

lemma *ex-det-eps*:
assumes $0 < e$
shows $\exists d \in D_D. \mathcal{L}_b v \leq L (mk\text{-dec-det } d) v + e *_R 1$
 $\langle \text{proof} \rangle$

lemma *thm-6-2-11*:
assumes $\text{eps} > 0$
shows $\exists d \in D_D. \nu_b\text{-opt} \leq \nu_b (mk\text{-stationary-det } d) + \text{eps} *_R 1$
 $\langle \text{proof} \rangle$

lemma *ex-det-dist-eps*:
assumes $0 < (e :: \text{real})$
shows $\exists d \in D_D. \text{dist } (\mathcal{L}_b v) (L (mk\text{-dec-det } d) v) \leq e$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-opt-le-det}$: $\nu_b\text{-opt } s \leq (\bigsqcup_{d \in D_D}. \nu_b (mk\text{-stationary-det } d) s)$
 $\langle \text{proof} \rangle$

lemma $\nu_b\text{-opt-eq-det}$: $\nu_b\text{-opt } s = (\bigsqcup_{d \in D_D}. \nu_b (mk\text{-stationary-det } d) s)$
 $\langle \text{proof} \rangle$

lemma *lemma-6-3-1-a*:
assumes $v0 \in \text{bfun}$
shows *uniform-limit UNIV* $(\lambda n. ((\lambda v. \mathcal{L} (\text{Bfun } v)) \widetilde{\sim} n) v0) \nu\text{-opt}$
sequentially
 $\langle \text{proof} \rangle$

lemma *thm-6-3-1-b-aux*: $(\lambda n. \text{dist } ((\mathcal{L}_b \widetilde{\sim} n) v) ((\mathcal{L}_b \widetilde{\sim} (\text{Suc } n)) v))$
 $\longrightarrow 0$
 $\langle \text{proof} \rangle$

definition *max-L-ex* $s v \equiv \text{has-arg-max } (\lambda a. L_a a v s) (A s)$

end

6.15 More Restrictive MDP Locales

locale *MDP-att- \mathcal{L}* = *MDP-reward* $A K r l$
for
 A **and**
 $K :: 's :: \text{countable} \times 'a :: \text{countable} \Rightarrow 's \text{ pmf}$ **and**
 r **and** $l +$
assumes *Sup-att: max-L-ex* $(s :: 's) v$
begin

theorem *thm-6-2-10-a-aux'*:
fixes $v :: 's \Rightarrow_b \text{real}$
assumes *is-arg-max* $(\lambda a. L_a a v s) (\lambda a. a \in A s) a$
shows $\mathcal{L}_b v s = L_a a v s$
 $\langle \text{proof} \rangle$

end

locale *MDP-act* = *MDP-att- \mathcal{L}* $A K r l$
for $A :: 's :: \text{countable} \Rightarrow ('a :: \text{countable}) \text{ set}$ **and** $K r l +$
fixes *arb-act* $:: 'a \text{ set} \Rightarrow 'a$
assumes *arb-act-in[simp]*: $X \neq \{\} \implies \text{arb-act } X \in X$
begin

definition *is-opt-act* $v s \equiv \text{is-arg-max } (\lambda a. L_a a v s) (\lambda a. a \in A s)$
abbreviation *opt-acts* $v s \equiv \{a. \text{is-opt-act } v s a\}$

lemma *is-opt-act-some*: *is-opt-act* $v s (\text{arb-act } (\text{opt-acts } v s))$
 $\langle \text{proof} \rangle$

lemma *some-opt-acts-in-A*: *arb-act* $(\text{opt-acts } v s) \in A s$
 $\langle \text{proof} \rangle$

lemma *ν -improving-opt-acts*: *ν -improving* $v0$ (*mk-dec-det* $(\lambda s. \text{arb-act}$

(*opt-acts* (*apply-bfun* *v0*) *s*))
 ⟨*proof*⟩

end

locale *MDP-finite-type* = *MDP-reward* *A K r l*
for *A* **and** *K* :: '*s*' :: *finite* × '*a*' :: *finite* ⇒ '*s*' *pmf* **and** *r l*

context *MDP-reward*
begin

lemma *ν_b-fin-zero[simp]*: *ν_b-fin* *p 0* = 0
 ⟨*proof*⟩

lemma *ν_b-fin-Suc[simp]*: *ν_b-fin* (*mk-stationary* *d*) (*Suc* *n*) = *ν_b-fin*
 (*mk-stationary* *d*) *n* + ((*l* *_R *P*₁ *d*) $\overset{\sim}{\sim}$ *n*) (*r-dec_b* *d*)
 ⟨*proof*⟩

lemma *ν_b-fin-eq*: *ν_b-fin* (*mk-stationary* *d*) *n* = (∑ *i* < *n*. ((*l* *_R *P*₁
d) $\overset{\sim}{\sim}$ *i*)) (*r-dec_b* *d*)
 ⟨*proof*⟩

lemma *L-iter*: (*L* *d* $\overset{\sim}{\sim}$ *m*) *v* = *ν_b-fin* (*mk-stationary* *d*) *m* + ((*l* *_R
*P*₁ *d*) $\overset{\sim}{\sim}$ *m*) *v*
 ⟨*proof*⟩

lemma *bounded-stationary-ν_b-fin*: *bounded* ((*λx*. (*ν_b-fin* (*mk-stationary*
x) *N*) *s*) ' *X*)
 ⟨*proof*⟩

lemma *bounded-disc-P₁*: *bounded* ((*λx*. (((*l* *_R *P*₁ *x*) $\overset{\sim}{\sim}$ *m*) *v*) *s*) ' *X*)
 ⟨*proof*⟩

lemma *bounded-disc-P₁'*: *bounded* ((*λx*. ((*P*₁ *x* $\overset{\sim}{\sim}$ *m*) *v*) *s*) ' *X*)
 ⟨*proof*⟩

lemma *L-iter-le-L_b*: *is-dec* *d* ⇒ (*L* *d* $\overset{\sim}{\sim}$ *n*) *v* ≤ (*L_b* $\overset{\sim}{\sim}$ *n*) *v*
 ⟨*proof*⟩

end

context *MDP-att-L*
begin

lemma *L_a-le-arg-max*: *a* ∈ *A* *s* ⇒ *L_a* *a v s* ≤ *L_a* (*arg-max-on* (*λa*.
L_a *a v s*) (*A* *s*)) *v s*
 ⟨*proof*⟩

lemma *arg-max-on-in*: $\text{has-arg-max } f \ Q \implies \text{arg-max-on } f \ Q \in Q$
<proof>

lemma *\mathcal{L}_b -eq- L_a -max*: $\mathcal{L}_b \ v \ s = L_a \ (\text{arg-max-on } (\lambda a. L_a \ a \ v \ s) \ (A \ s))$
v s
<proof>

lemma *ex-opt-det*: $\exists d \in D_D. \mathcal{L}_b \ v = L \ (\text{mk-dec-det } d) \ v$
<proof>

lemma *ex-improving-det*: $\exists d \in D_D. \nu\text{-improving } v \ (\text{mk-dec-det } d)$
<proof>

end

end

References

- [1] J. Hölzl and T. Nipkow. Markov models. *Archive of Formal Proofs*, Jan. 2012. https://isa-afp.org/entries/Markov_Models.html, Formal proof development.
- [2] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley, 1994.