

Supplementary Materials

The following content was not necessarily subject to peer review.

8 Proofs of Technical Analysis

In this section, we provide formal proofs to our technical analysis in detail. For better legibility, we first recall the equations and results that we need for our proofs.

$$\forall \omega \in \Omega, \quad \max_{\pi \in \Pi} \phi(\mathbf{V}(\pi)), \quad (3)$$

where Ω is the set of valid preference weights sorted in descending order, $\mathbf{V}(\pi) = \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t]$ is the expected discounted return, and $\phi(\mathbf{J}) = \sum_{i=1}^N w_i \mathbf{V}_{(i)}$ with $\mathbf{V}_{(1)} \leq \dots \leq \mathbf{V}_{(n)}$.

Lemma 8.1. *For any MOMDP with linear preferences over objectives, the CCS contains an optimal policy for any linear combination of the objectives.*

Proof. Let \mathcal{S} be the state space, \mathcal{A} be the action space, and $\mathbf{r} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbf{r}^N$ be the vector-valued reward function, where N is the number of objectives. Consider a linear preference vector $\omega \in \Omega$, where $\Omega = \{\omega \in \mathbf{r}^N : \sum_{i=1}^N w_i = 1, w_i \geq 0\}$. For any policy π , the expected return under a preference ω is given by $\omega(\mathbb{E}_\pi[\sum_{t=1}^{\infty} \gamma^{t-1} \mathbf{r}(s_t, a_t) \mid s_0 = s])$. Thus, the optimal policy π_ω^* for preference ω satisfies

$$\pi_\omega^* = \operatorname{argmax}_{\pi} \omega^T \mathbf{V}^\pi(s), \quad \forall s \in \mathcal{S}.$$

By the definition of the CCS, for any $\omega \in \Omega$, there exists a policy $\pi_{\text{CCS}} \in \text{CCS}$ such that

$$\omega^T \mathbf{V}^{\pi_{\text{CCS}}}(s) \geq \omega^T \mathbf{V}^\pi(s), \quad \forall \pi \in \Pi, \forall s \in \mathcal{S}.$$

To prove the proposition, let's recall the Convex Hull Value Iteration (CHVI) algorithm (Barrett & Narayanan, 2008). Note that the CHVI algorithm iteratively updates the value function for each state by considering the convex hull of the achievable rewards via

$$\mathbf{V}(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathbf{P}(s' \mid s, a) \text{CH}(\mathbf{r}(s, a) + \gamma \mathbf{V}(s')),$$

where $\text{CH}(\cdot)$ denotes the convex hull operation. This update rule ensures that the value function $\mathbf{V}(s)$ lies within the convex hull of the achievable rewards and the $\text{CH}(\cdot)$ achievable value functions $\mathbf{V}^\pi(s) \mid \pi \in \Pi$ forms the CCS. Therefore, for any linear preference vector ω , there must exist at least a policy π_{CCS} such that

$$\omega^T \mathbf{V}^{\pi_{\text{CCS}}}(s) = \max_{\pi \in \Pi} \omega^T \mathbf{V}^\pi(s), \quad \forall s \in \mathcal{S}.$$

The resulting policies form the CSS, which are sufficient to cover all linear preferences $\omega \in \Omega$. Thus, for any linear combination of objectives, the optimal policy can be found within the CSS, confirming its sufficiency and optimality. \square

While GGF introduces non-linear fairness objectives, its piecewise linearity and concavity allow representation as a maximum of linear functions, which ensures that solutions lie within the CCS. The following proposition establishes the sufficiency of the CCS in representing optimal policies for ϕ_{GGF} preference weights.

Proposition 8.1. *For any $s \in \mathcal{S}$ in an MOMDP and a piecewise-linear concave welfare function ϕ_{GGF} (e.g., GGF) that can be represented as, $\phi_{\text{GGF}}(\mathbf{V}^\pi(s)) = \min_{\sigma \in \mathbb{S}_N} \{\omega_\sigma^\top \mathbf{V}^\pi(s)\}$, there exists a policy $\pi^* \in \text{CCS}$ such that:*

$$\phi_{\text{GGF}}(\mathbf{V}^{\pi^*}(s)) \geq \phi_{\text{GGF}}(\mathbf{V}^\pi(s)) \quad \forall \pi \in \Pi.$$

636 *Proof.* Policies $\pi^* \in CCS$ satisfies

$$\forall \omega, \quad \exists \pi^* \in CCS, \quad \forall \pi \in \Pi : \omega^\top \mathbf{V}^{\pi^*}(s) \geq \omega^\top \mathbf{V}^\pi(s).$$

637 GGF can be represented as

$$\phi_{\text{GGF}}(\mathbf{V}^\pi(s)) = \min_{\sigma \in \mathbb{S}_N} \{\omega_\sigma^\top \mathbf{V}^\pi(s)\}.$$

638 and that the set \mathbb{S}_N consists of the vectors $\{\omega_1, \omega_2, \dots, \omega_{N!}\}$. Therefore, for any policy π we can
639 write

$$\phi_{\text{GGF}}(\mathbf{V}^\pi(s)) = \min \{\omega_1^\top \mathbf{V}^\pi(s), \omega_2^\top \mathbf{V}^\pi(s), \dots, \omega_{N!}^\top \mathbf{V}^\pi(s)\}.$$

640 Given $\pi^* \in CCS$, $\forall \pi \in \Pi$ satisfies

$$\omega_i^\top \mathbf{V}^{\pi^*}(s) \geq \omega_i^\top \mathbf{V}^\pi(s), \text{ for all } i = 1, 2, \dots, N!,$$

641 it immediately follows that

$$\min_{1 \leq i \leq N!} \omega_i^\top \mathbf{V}^{\pi^*}(s) \geq \min_{1 \leq i \leq N!} \omega_i^\top \mathbf{V}^\pi(s).$$

642 Thus, we conclude that,

$$\forall \pi \in \Pi, \quad \exists \pi^* \in CCS, \quad \phi_{\text{GGF}}(\mathbf{V}^{\pi^*}(s)) \geq \phi_{\text{GGF}}(\mathbf{V}^\pi(s)).$$

643

□

644 **Fairness of Non-Stationary Policies.** In fair MORL, learning non-stationary policies can be par-
645 ticularly beneficial, as they leverage historical information to make more informed decisions and
646 adapt over time (see Section 4).

647 **Proposition 8.2.** Let the reward \mathbf{r} be nonnegative, and Π_S and Π_{NS} be the sets of stationary and
648 non-stationary policies, respectively. For any $s \in \mathcal{S}$ in an MOMDP and a given ϕ_{GGF} , there exists a
649 non-stationary policy $\pi_{NS} \in \Pi_{NS}$ that achieves a higher welfare score than any stationary policy
650 $\pi_S \in \Pi_S$, i.e.,

$$\exists \pi_{NS} \in \Pi_{NS} : \phi_{\text{GGF}}(\mathbf{V}^{\pi_{NS}}(s)) \geq \max_{\pi_S \in \Pi_S} \phi_{\text{GGF}}(\mathbf{V}^{\pi_S}(s))$$

651 *Proof.* Let the state value function be defined by:

$$\mathbf{V}(s) = \mathbb{E} [\mathbf{G}_t | s_t = s]$$

652 where the return \mathbf{G}_t is given by:

$$\mathbf{G}_t = \sum_{k=0}^{\infty} \gamma^k \mathbf{r}_{t+k+1}.$$

653 Suppose an episode begins at time t and terminates at time T_{end} . For any intermediate time T with
654 $t \leq T < T_{\text{end}}$, we can decompose the return into two parts:

$$\mathbf{G}_t = \underbrace{\mathbf{r}_{t+1} + \gamma \mathbf{r}_{t+2} + \dots + \gamma^{T-t-1} \mathbf{r}_T}_{\mathbf{G}_t^{(1)}} + \underbrace{\gamma^{T-t} (\mathbf{r}_{T+1} + \gamma \mathbf{r}_{T+2} + \dots)}_{\mathbf{G}_t^{(2)}}.$$

655 With above decomposition, We define value function as two parts:

$$\text{Early-period value function: } \mathbf{V}_1(s) = \mathbb{E} [\mathbf{G}_t^{(1)} | s_t = s]$$

656

$$\text{Late-period value function: } V_2(s) = \mathbb{E} \left[G_t^{(2)} \middle| s_T = s \right]$$

657 so that

$$V(s) = V_1(s) + \gamma^{T-t} V_2(s)$$

658 At time T , stationary policy π_S selects action solely based on late period value function $V_2(s)$, while
 659 non-stationary policy has access to both early $V_1(s)$ and late period value function $V_2(s)$ and can
 660 condition its action selection on the combined information given by two value functions.

661 Under a stationary policy, The total value can be presented as:

$$V^{\pi_S}(s) = V_1(s) + \gamma^{T-t} \underset{V_2(s)}{\operatorname{argmax}} \{ \phi_{\text{GGF}}[V_2(s)] \}$$

662 In contrast, under a non-stationary policy the total value is given by

$$V^{\pi_{NS}}(s) = \underset{V_1(s), V_2(s)}{\operatorname{argmax}} \{ \phi_{\text{GGF}}[V_1(s) + \gamma^{T-t} V_2(s)] \}$$

663 therefore:

$$\exists \pi_{NS} \in \Pi_{NS} : \phi_{\text{GGF}}(V^{\pi_{NS}}(s)) \geq \max_{\pi_S \in \Pi_S} \phi_{\text{GGF}}(V^{\pi_S}(s))$$

664 This completes the proof. \square

665 **Optimality of Stochastic Policies for Fairness** Unlike the single-objective scenario, in MORL, a
 666 deterministic policy may not be optimal. A fairer solution can often be achieved through random-
 667 ization.

668 **Proposition 8.3.** *Let Π_{ST} be the set of stochastic policies and Π_D be the set of deterministic policies.*
 669 *For an MOMDP \mathcal{M} and a concave welfare function such as ϕ_{GGF} , there exists a stochastic policy*
 670 *$\pi_{ST} \in \Pi_{ST}$ such that:*

$$\phi_{\text{GGF}}(V^{\pi_{ST}}) \geq \max_{\pi_D \in \Pi_D} \phi_{\text{GGF}}(V^{\pi_D}).$$

671 *Proof.* The key idea here is that a stochastic policy can represent a convex combination of determin-
 672 istic policies for any concave welfare function ϕ_{GGF} [Busa-Fekete et al. \(2017\)](#). Hence, stochastic
 673 policies can achieve outcomes in the objective space that are unattainable by deterministic policies.
 674 Specifically, for ϕ_{GGF} , a deterministic policy π_D yields a fixed utility vector V^{π_D} while a stochastic
 675 policy π_{ST} can yield a distribution over utility vectors. Thanks to concavity of ϕ_{GGF} , which makes
 676 our problem in 2 convex optimization and Jensen's inequality ([Jensen, 1967](#)), we obtain

$$\phi_{\text{GGF}}(\mathbb{E}_{\tau \sim \pi}[V^{\pi_{st}}]) \geq \mathbb{E}_{\tau \sim \pi}[\phi_{\text{GGF}}(V^{\pi_{st}})]. \quad (4)$$

677 Since ϕ_{GGF} is a piecewise linear concave function, there exists a stochastic policy π_{st} that is a convex
 678 combination of deterministic policies such that

$$\mathbb{E}_{\tau \sim \pi}[\phi(V^{\pi_{st}})] \geq \max_{\pi_d \in \Pi_D} \phi(V^{\pi_d}). \quad (5)$$

679 By combining (4) and (5), we can obtain

$$\phi(\mathbb{E}_{\tau \sim \pi}[V^{\pi_{st}}]) \geq \mathbb{E}_{\tau \sim \pi}[\phi(V^{\pi_{st}})] \geq \max_{\pi_d \in \Pi_D} \phi(V^{\pi_d}).$$

680 This completes the proof. \square

681 The optimality of stochastic policies implies that restricting the search for fair solutions to deter-
 682 ministic policies is insufficient. Stochastic policies offer a broader range of solutions and may better
 683 capture the trade-offs among multiple objectives, enhancing the overall fairness of the policy.

9 Fairness

In a fair single-policy setting, where the goal is to learn a single policy treating all users equally, three fairness principles, efficiency, equity, and impartiality, are defined below.

Definition 9.1. *Efficiency states that among two feasible solutions, if one solution is (weakly or strictly) preferred by all users, then it should be preferred to the other one, e.g., $\mathbf{u} \succ \mathbf{u}' \Rightarrow \phi(\mathbf{u}) > \phi(\mathbf{u}')$, where $\phi(\mathbf{u})$ is the scalar utility function that specifies the value of a solution.*

Intuitively, the efficiency property specifies that given all else equal, one prefers to increase a user's utility. In the MORL setting, the efficiency property simply means Pareto dominance. More specifically, a solution is considered efficient if it is not dominated by any other solution for all objectives.

Next, we discuss the significance of the *equity* property, which is a stronger property than efficiency and is often associated with distributive justice, as it refers to the fair distribution of resources or opportunities. This property ensures that a fair solution follows the *Pigou-Dalton principle* (Moulin, 2004), which states the transferring of rewards from the more advantaged users to the less advantaged users.

Definition 9.2. *A solution satisfies the Pigou-Dalton principle if for all \mathbf{u}, \mathbf{u}' equal except for $u_i = u'_i + \delta$ and $u_j = u'_j - \delta$ where $u'_i - u'_j > \delta > 0$, $\phi(\mathbf{u}) > \phi(\mathbf{u}')$.*

Finally, we discuss the *impartiality* property. This property is rooted in the principle of “equal treatment of equals”, which states that individuals sharing similar characteristics should be treated similarly.

Definition 9.3. *In a system, individuals with similar characteristics should be treated similarly, i.e., the solution should be independent of the order of its arguments $\phi(\mathbf{u}) = \phi(\mathbf{u}_\sigma)$, where σ is a permutation and \mathbf{u}_σ is the vector obtained from vector \mathbf{u} permuted by σ .*

9.1 Welfare Function

A welfare function, denoted as $\phi : \mathbb{R}^D \rightarrow \mathbb{R}$, aggregates the utilities of all users (or objectives) and offers a metric of the overall desirability of a solution for the entire group, where ω represents the set of aggregation weights for all objectives. One well-established welfare function used in this paper is the generalized Gini welfare function. The generalized Gini welfare function constitutes a specific instance of the ordered weighted average (OWA) (Yager, 1988). It is a renowned welfare function employed in multi-objective optimization (Weng, 2019; Siddique et al., 2020; Zimmer et al., 2021; Do & Usunier, 2022; Yu et al., 2023a;b; Siddique et al., 2023), initially devised to quantify income distribution inequality in economics (Weymark, 1981). The generalized Gini welfare function is defined as follows:

$$G_\omega(\mathbf{u}) = \sum_{i=1}^N \omega_{\sigma(i)} u_i = \mathbf{w}_\sigma^T \mathbf{u}, \quad (6)$$

where $\sigma \in \mathbb{S}_N$, which depends on ω , is the permutation that sorts the components of ω and $\omega_\sigma = (\omega_{\sigma(1)}, \dots, \omega_{\sigma(N)})$. Equation (6) holds as the weights are rearranged based on the utility vector, assigning the largest weight to the smallest component of \mathbf{u} , the second-largest weight to the second-smallest component of \mathbf{u} , and so forth.

The generalized Gini welfare function satisfies the three fairness properties. Due to the positive weights, it is monotonically related to Pareto dominance, fulfilling the efficiency property. Moreover, the reordering of the components in the welfare function makes it symmetric with respect to its components, satisfying the impartiality property. Lastly, as the generalized Gini weights are positive and decreasing, it is Schur-concave, meeting the equity property.

Among numerous welfare functions, the generalized Gini welfare function possesses several favorable properties, namely, simplicity as it is a weighted sum in the Lorenz space (Chakravarty, 1990;

Perny et al., 2013), well-understood properties axiomatized by Weymark (1981), and generality. These favorable properties make it a suitable choice for addressing the challenge of finding fair solutions. Moreover, it is notably a concave function, which will make the solution to our problem easier.

To emphasize the versatility of the generalized Gini welfare function, various special cases can be derived by adjusting its weights accordingly. These cases include:

- **Maxmin fairness:** Setting $\omega_1 = 1$ and $\omega_i = 0$ for $i = 2, \dots, K$ corresponds to the maxmin notion of fairness (Rawls, 1971).
- **Regularized maxmin fairness:** Assigning $\omega_1 = 1$ and $\omega_i = \varepsilon$ for $i = 2, \dots, K$ aligns with the regularized maxmin notion of fairness.
- **Utilitarian approach:** Setting $\omega_i = 1/K$ represents the utilitarian approach.
- **Leximin fairness:** If the ratio ω_j/ω_{j+1} tends toward infinity, it corresponds to the leximin notion of fairness (Rawls, 1971; Kurokawa et al., 2015).

10 Descriptions of Environments

10.1 Species Conservation

In the field of ecology, the challenge of conserving interdependent endangered species is paramount. The simulation environment focuses on the balance required in the conservation of two such species: the sea otter and the northern abalone, which are currently endangered. The predation relationship between these species, with sea otters feeding on abalones, presents a unique challenge that requires careful consideration of fairness and equity in conservation efforts. Based on the framework in (Chadès et al., 2012), we define the state space as the current population numbers of the sea otters and northern abalones. The action space consists of: introducing sea otters, enforcing antipoaching measures, controlling sea otter populations, implementing a combination of half-antipoaching and half-controlled sea otters, or taking no action. Each action carries significant ecological consequences; for instance, while the reintroduction of sea otters is essential for maintaining the abalone population, it must be carefully managed to prevent the abalone’s extinction. Conversely, overlooking other management actions could lead to the demise of either species. The transition function employed in our model accounts for population dynamics, including external threats such as poaching and oil spills. Since our objective is to optimize the population densities of both species, we define the reward function as the densities of both species, i.e., $N = 2$.

10.2 Resource Gathering

In this scenario of resource gathering, we consider a 5×5 grid world domain inspired from (Barrett & Narayanan, 2008). This domain presents a unique challenge centered around the acquisition of three types of resources: gold, gems, and stones, thereby establishing a multi-objective framework with $\mathcal{K} = 3$. The autonomous agent is positioned within this grid world, and resources are distributed randomly across various locations. As a resource is collected by the agent, it is immediately regenerated at a new random location within the grid, ensuring a perpetual availability of resources. In this problem, the state is characterized by the agent’s current location on the grid and a cumulative count of each type of resource collected over the course of the agent’s trajectory. The agent can navigate the grid through actions aligned with the four cardinal directions: up, down, left, and right, facilitating movement across the grid. To add complexity to the resource management challenge, resources are assigned differing values, reflecting their relative importance. Specifically, gold and gems are attributed a value of 1, underscoring their significance, whereas stones are considered less valuable, with a value of 0.4. This valuation leads to an intentionally uneven distribution of resources within the grid, comprising two stones, one gold, and one gem. This configuration is designed to simulate a scenario where the agent must not only maximize the collection of resources but also achieve a balanced acquisition across the different types of resources. The overarching objective for the agent in this environment is dual: to maximize the total value of resources collected while

ensuring an equitable collection across the various resource types. Achieving this balance is crucial for optimizing the agent’s resource-gathering strategy, enhancing its overall utility and adaptability within the dynamic grid world. This nuanced approach to resource management in a simulated environment offers insights into the complexities of resource distribution and acquisition strategies, contributing to the broader discourse on multi-objective optimization in dynamic settings.

10.3 Multi-Product Web Advertising

We now consider the multi-product web advertising (MWP) problem, where an online store offers N distinct types of products for sale and an intelligent agent makes strategic decisions at each timestep about which advertisement to display: a product-specific advertisement for one of the products $i \in [0, \dots, N - 1]$, or a general advertisement that is not tailored to any specific product. The effectiveness of an advertisement is contingent upon its relevance to the customer’s recent web activity, with appropriate advertisements significantly increasing the likelihood of a purchase, whereas inappropriate ones may deter the customer altogether. The state space of this problem is defined by the number of products available in the store, augmented by the number of visits, purchases, and exits. A visit state indicates a customer’s interest in a particular product, a purchase state signifies the completion of a transaction, and an exit state occurs when a customer leaves the website without making a purchase. The action space is expanded to $n + 1$ actions, where actions 0 through n correspond to displaying advertisements for specific products, and action n represents the option to show a general advertisement that does not target any specific product in the inventory. This additional action introduces an additional layer of complexity, as the agent must decide the optimal moment to transition between states. The reward function is designed such that the agent receives a reward of 1 in the i^{th} dimension of the reward vector if a product of type i is sold after the display of its advertisement. The primary objective of this problem is to maximize the aggregate returns from product sales while striving for an equitable distribution of sales across the different product types. This goal underscores the need for fair solutions that not only optimize overall profitability but also ensure a balanced representation of product sales, thereby addressing the dual challenges of efficiency and equity in this domain.

11 Hyperparameters

To ensure reproducibility, we have meticulously documented all hyperparameters across different environments in Tables 1, 2, 3, and 4. We utilize the well-known high-quality MORL baselines¹ for implementing baseline algorithms. In these tables, we present the hyperparameters corresponding to Envelope, GPI, PCN, and our proposed algorithms in three distinct environments, namely, species conservation (SC), resource gathering (RC), and multi-web product advertising (MWP).

¹<https://github.com/LucasAlegre/morl-baselines>

Table 1: Set of hyperparameters used for training Envelope.

Hyperparameter	SC	RC	MWP
Discount factor (γ)	0.99	0.99	0.99
Learning rate (α)	0.0001	0.0005	0.005
Batch size	64	64	64
Hidden Layers	256 x 256 x 256 x 256	256 x 256 x 256 x 256	256 x 256 x 256 x 256
Buffer Size	50000	50000	50000
Initial Epsilon	1.0	1.0	1.0
Final Epsilon	0.05	0.05	0.05
Epsilon Decay Steps	50000	50000	50000
Learning Starts	100	100	100
Gradient Updates	1	1	5
Max Gradient Norm	1.0	1.0	1.0
Ω Distribution	Gaussian	Gaussian	Gaussian
Tau	0.5	0.5	0.5

Table 2: Set of hyperparameters used for training our proposed methods.

Hyperparameter	SC	RC	MWP
Discount factor (γ)	0.99	0.99	0.99
Learning rate (α)	0.0001	0.0005	0.005
Batch size	64	64	64
Hidden Layers	256 x 256 x 256 x 256	256 x 256 x 256 x 256	256 x 256 x 256 x 256
Buffer Size	50000	50000	50000
Initial Epsilon	1.0	1.0	1.0
Final Epsilon	0.05	0.05	0.05
Epsilon Decay Steps	50000	50000	50000
Learning Starts	100	100	100
Gradient Updates	1	1	5
Max Gradient Norm	1.0	1.0	1.0
Ω Distribution	Gaussian	Gaussian	Gaussian
Tau	0.5	0.5	0.5

Table 3: Set of hyperparameters used for training GPI.

Hyperparameter	SC	RC	MWP
Discount factor (γ)	0.99	0.99	0.99
Learning rate (α)	0.0001	0.0005	0.005
Batch size	128	128	256
Hidden Layers	256 x 256 x 256 x 256	256 x 256 x 256 x 256	256 x 256 x 256 x 256
Num Networks	2	2	2
Buffer Size	50000	50000	50000
Initial Epsilon	1.0	1.0	1.0
Final Epsilon	0.05	0.05	0.05
Epsilon Decay Steps	50000	50000	50000
Learning Starts	100	100	100
Gradient Updates	1	1	5

Table 4: Set of hyperparameters used for training PCN.

Hyperparameter	SC	RC	MWP
Discount factor (γ)	0.99	0.99	0.99
Learning rate (α)	0.0001	0.0001	0.0005
Batch size	128	256	128
Hidden Layers	64 x 64	64 x 64	64 x 64
Desired Return	[1, 1]	[200, 200, 200]	[100, 100, 100, 100, 100]
Buffer Size	500000	500000	1000000
Max Horizon	5000	1000	1000