# <span id="page-0-2"></span>Supplemental material for the paper: Constrained Cross-Entropy Method for Safe Reinforcement Learning

We repeat some notations and definitions here and restate Problem 1.

For  $\rho \in (0,1)$ ,  $v \in V$  and any function  $X : \Theta \to \mathbb{R}$ , we denote the  $\rho$ -quantile of X for  $\theta \sim f_v$  by  $\xi_X(\rho, \nu)$ . We also define  $\delta : \mathbb{R} \times \{\geq, \leq, >, <, =\} \times \mathbb{R} \to \{0, 1\}$  as an indicator function such that for  $\circ \in \{\geq, \leq, >, <, =\}, \delta(x \circ y) = 1$  if and only if  $x \circ y$  holds. The surrogate objective function for the unconstrained CE method is  $\mathbb{E}_{\theta \sim f_{\boldsymbol{v}}(\cdot)}[G(\pi_{\theta})\delta(G(\pi_{\theta}) \geq \xi_G(1-\rho,\boldsymbol{v}))]$ . In other words, a policy  $\pi_{\theta}$  is considered as highly ranked if  $G(\pi_{\theta}) \ge \xi_G(1-\rho, v)$ . When there is a constraint  $H(\pi) \le d$ , we define  $U : \Pi_{\Theta} \to \mathbb{R}$  such that  $U(\pi_{\theta}) := G(\pi_{\theta})\delta(H(\pi_{\theta}) \leq d)$  for any  $\theta \in \Theta$  and extend the surrogate function as follows:

<span id="page-0-0"></span>
$$
L(\boldsymbol{v};\rho) := \begin{cases} \mathbb{E}_{\theta \sim f_{\boldsymbol{v}}(\cdot)} [G(\pi_{\theta})\delta(H(\pi_{\theta}) \leq \xi_H(\rho, \boldsymbol{v}))], & \text{if } \xi_H(\rho, \boldsymbol{v}) > d; \\ \mathbb{E}_{\theta \sim f_{\boldsymbol{v}}(\cdot)} [U(\pi_{\theta})\delta(U(\pi_{\theta}) \geq \xi_U(1-\rho, \boldsymbol{v}))], & \text{otherwise.} \end{cases}
$$
(1)

We can combine the two cases. Define  $S : \Pi_{\Theta} \times \mathcal{V} \times (0, 1) \rightarrow \{0, 1\}$  such that

$$
S(\pi_{\theta}, \mathbf{v}, \rho) := \delta(\xi_H(\rho, \mathbf{v}) > d) \delta(H(\pi_{\theta}) \leq \xi_H(\rho, \mathbf{v})) +
$$
  

$$
\delta(\xi_H(\rho, \mathbf{v}) \leq d) \delta(H(\pi_{\theta}) \leq d) \delta(U(\pi_{\theta}) \geq \xi_U(1 - \rho, \mathbf{v})),
$$

then [\(1\)](#page-0-0) can be rewritten as

<span id="page-0-1"></span>
$$
L(\boldsymbol{v};\rho) = \mathbb{E}_{\theta \sim f_{\boldsymbol{v}}(\cdot)}[G(\pi_{\theta})S(\pi_{\theta}, \boldsymbol{v}, \rho)].
$$
\n(2)

We first restate Problem 1 here.

**Problem 1.** *Given a set*  $\Pi = {\pi_\theta : \theta \in \Theta}$  *of policies with parameter space*  $\Theta$ *, an NEF*  $F_V$  =  ${f_{\boldsymbol{v}}(\cdot) \in \mathcal{D}(\Theta) : \boldsymbol{v} \in \mathcal{V}}$  *of distributions over*  $\Theta$ , two functions  $G: \Pi \to \mathbb{R}^+$  and  $H: \Pi \to \mathbb{R}$ , a *constraint upper bound*  $\hat{d}$  *and*  $\rho \in (0,1)$ *, compute*  $v^* \in V$  *such that* 

$$
\boldsymbol{v}^* = \arg\max_{\boldsymbol{v}\in\mathcal{V}} L(\boldsymbol{v};\rho),
$$

*where*  $L: V \times (0,1) \rightarrow \mathbb{R}$  *is defined in* [\(2\)](#page-0-1).

<span id="page-0-3"></span>*Remark* 1. For technical reasons, we approximate the binary function  $\delta$  with a Lipschitz continuous piecewise linear function  $\tilde{\delta}_{\varepsilon} : \mathbb{R} \times \{ \geq, \leq, >, <, = \} \times \mathbb{R} \to [0, 1]$  with  $\varepsilon > 0$ . For example, for any  $(x, y) \in \mathbb{R}^2$ :

$$
\tilde{\delta}_{\varepsilon}(x \ge y) = \begin{cases} 1, & \text{if } x \ge y \\ \frac{x-y}{\varepsilon}, & \text{if } y > x \ge y - \varepsilon \\ 0, & \text{otherwise.} \end{cases}
$$

With slight abuse of notation, we use  $\delta$  to represent  $\tilde{\delta}_{\varepsilon}$  for some small enough  $\varepsilon$ . Therefore  $\delta(x \circ y)$ is Lipschitz continuous both in  $x$  and in  $y$ .

# Proofs for Section 4.2

In this section we provide proofs for all theorems in our paper. The main idea behind the proof of Theorem 4.1 is similar to that of Theorem 3.1 [Hu et al.](#page-10-0) [\[2012\]](#page-10-0), although the details are adapted to our problem.

32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montreal, Canada. ´

The following lemma shows that Assumption [1](#page-0-2) in Section [4.2](#page-0-2) is sufficient to guarantee that  $m^{-1}$ exists and is continuously differentiable.

<span id="page-1-0"></span>**Lemma 1.** Let  $F_v$  be an NEF. Define  $m : \mathbb{R}^{d_v} \to \mathbb{R}^{d_v}$  such that for each  $v \in V := \{v \in V\}$  $\mathbb{R}^{d_{\bm v}}:K(\bm v)<\infty\},\,m(\bm v):=\mathbb{E}_{\bm v}[\Gamma(\theta)].$  Then  $m^{-1}$  exists and is continuously differentiable over  $\{\eta : \exists v \in int(V) \text{ s.t. } \eta = m(v)\}.$ 

*Proof.*

$$
\frac{\partial}{\partial v}m(v) = \frac{\partial}{\partial v} \int_{\Theta} \Gamma(\theta) f_{v}(\theta) d\theta
$$
\n
$$
= \int_{\Theta} \left(\frac{\partial}{\partial v} f_{v}(\theta)\right) (\Gamma(\theta))^{\mathsf{T}} d\theta
$$
\n
$$
= \int_{\Theta} \left(f_{v}(\theta)(\Gamma(\theta) - m(v))(\Gamma(\theta))^{\mathsf{T}} d\theta\right)
$$
\n
$$
= \int_{\Theta} f_{v}(\theta) \Gamma(\theta) \Gamma(\theta) \Gamma(d\theta - m(v) \cdot \left(\int_{\Theta} f_{v}(\theta) \Gamma(v) dv\right)^{\mathsf{T}}
$$
\n
$$
= \mathbb{E}_{v} [(\Gamma(\theta) - m(v)) (\Gamma(\theta) - m(v))^{\mathsf{T}}]
$$
\n
$$
= \text{Cov}_{v} [\Gamma(\theta)],
$$

which is positive definite by Assumption [1.](#page-0-2) As both  $f_v(\cdot)$  and  $m(v)$  are continuous at  $v \in int(V)$ ,  $m(v)$  is continuously differentiable with nonzero derivative at any  $v \in int(V)$ . Therefore by inverse mapping theorem,  $m^{-1}$  exists and is continuously differentiable over  $\{\eta : \exists v \in int(V) s.t. \eta =$  $m(\boldsymbol{v})\}.$ 

The following lemma derives [\(5\)](#page-0-2) in Section [4.2.](#page-0-2)

**Lemma 2.** Let  $F_v$  be an NEF,  $v \in V$  and  $g \in \mathcal{D}(\Theta)$ . Let  $v' := \arg \min_{v \in V} D_{KL}(\tilde{g}_v, f_v)$ . Assume *that*  $v' \in int(V)$ *, then* 

$$
m(\boldsymbol{v}') - m(\boldsymbol{v}) = -\alpha \Big(\frac{\partial}{\partial \boldsymbol{v}''} D_{KL}(g_{\boldsymbol{v}}^*, f_{\boldsymbol{v}''}) \Big) \Big|_{\boldsymbol{v}'' = \boldsymbol{v}}.
$$

*Proof.* As  $F_v$  is an NEF,  $\log f_{\mathbf{v}''}(\theta) = (\mathbf{v}'')^T \Gamma(\theta) - K(\mathbf{v}'')$  is concave in  $\mathbf{v}''$  and thus  $-\int_{\Theta} \tilde{g}_{\boldsymbol{v}}(\theta) \log f_{\boldsymbol{v}''}(\theta) d\theta$  is convex in  $\boldsymbol{v}''$ . Therefore  $\boldsymbol{v}'$  can be found by setting the gradient to zero.

Let  $-\frac{\partial}{\partial v''}\int_{\Theta} \tilde{g}_{v}(\theta) \log f_{v''}(\theta) d\theta = 0$ . By Assumption [\(2d\),](#page-0-2)  $\Theta$  is bounded. By definition of NEF,  $f_{\mathbf{v}^{\prime\prime}}$  is continuously differentiable. Then by dominated convergent theorem,

$$
- \frac{\partial}{\partial v''} \int_{\Theta} \tilde{g}_{\boldsymbol{v}}(\theta) \log f_{\boldsymbol{v}''}(\theta) d\theta
$$
  
\n
$$
= - \int_{\Theta} \tilde{g}_{\boldsymbol{v}}(\theta) \frac{\partial}{\partial v''} \log f_{\boldsymbol{v}''}(\theta) d\theta
$$
  
\n
$$
= - \int_{\Theta} \tilde{g}_{\boldsymbol{v}}(\theta) (\Gamma(\theta) - m(\boldsymbol{v}'')) d\theta
$$
  
\n
$$
= - \int_{\Theta} (\alpha g_{\boldsymbol{v}}^*(\theta) + (1 - \alpha) f_{\boldsymbol{v}}(\theta)) \Gamma(\theta) d\theta + m(\boldsymbol{v}'')
$$
  
\n
$$
= - \alpha \Big( \int_{\Theta} g_{\boldsymbol{v}}^*(\theta) \Gamma(\theta) d\theta - m(\boldsymbol{v}) \Big) + m(\boldsymbol{v}'') - m(\boldsymbol{v})
$$
  
\n
$$
= - \alpha \Big( \mathbb{E}_{g_{\boldsymbol{v}}^*} [\Gamma(\theta)] - m(\boldsymbol{v}) \Big) + m(\boldsymbol{v}'') - m(\boldsymbol{v}),
$$

which equals to 0 when  $v'' = v'$ .

Since  $G(\pi_{\theta}) \ge 0$  for all  $\theta \in \Theta$ ,  $G(\pi_{\theta}) \equiv |G(\pi_{\theta})|$ . Therefore

$$
m(\mathbf{v}') = m(\mathbf{v}) + \alpha \Big( \int_{\Theta} g_{\mathbf{v}}^*(\theta) \Gamma(\theta) d\theta - m(\mathbf{v}) \Big) = (1 - \alpha) m(\mathbf{v}) + \alpha \mathbb{E}_{g_{\mathbf{v}}^*}[\Gamma(\theta)].
$$

On the other hand,

$$
\frac{\partial}{\partial v'} D_{KL}(g_v^*, f_{v'})
$$
\n
$$
= \frac{\partial}{\partial v'} \mathbb{E}_{g_v^*} \log \frac{g_v^*(\theta)}{f_{v'}(\theta)} |_{v'=v} = -\frac{\partial}{\partial v'} \mathbb{E}_{g_v^*} \log f_{v'}(\theta) |_{v'=v}
$$
\n
$$
= -\frac{\partial}{\partial v'} \mathbb{E}_{g_v^*} (\mathbf{v'}^{\mathsf{T}} \Gamma(\theta) - K(\mathbf{v'})) |_{v'=v}
$$
\n
$$
= -\frac{\partial}{\partial v'} (\mathbf{v'}^{\mathsf{T}} \mathbb{E}_{g_v^*} [\Gamma(\theta)]) |_{v'=v} + \frac{\partial}{\partial v'} K(\mathbf{v'}) |_{v'=v}
$$
\n
$$
= -\mathbb{E}_{g_v^*} [\Gamma(\theta)] + m(\mathbf{v}).
$$

Therefore

<span id="page-2-0"></span>
$$
m(\boldsymbol{v}') - m(\boldsymbol{v}) = -\alpha \Big(\frac{\partial}{\partial \boldsymbol{v}''} D_{KL}(g_{\boldsymbol{v}}^*, f_{\boldsymbol{v}''}) \Big) \Big|_{\boldsymbol{v}'' = \boldsymbol{v}}.
$$

 $\Box$ 

# Proof of Theorem [4.1](#page-0-2)

In practice we can only estimate expectations and quantiles using finite samples. Let  $\mathcal{Y}_l$  =  $\{\theta_1,\ldots,\theta_{n_l}\}$  be the set of samples in the  $l^{th}$  iteration with sampling distribution  $f_{v_l}$ . We denote the sample estimate of  $S(\pi_{\theta}, \nu, \rho)$  by  $\hat{S}(\pi_{\theta}, \nu, \rho)$ .

Consider the equation in the Step [11](#page-0-2) of Algorithm [1:](#page-0-2)

$$
\hat{\eta}_{l+1} = \alpha_l \frac{\sum_{i=1}^{n_l} G(\pi_{\theta_i}) \hat{S}(\pi_{\theta_i}, \boldsymbol{v}_l, \rho) \Gamma(\theta_i)}{\sum_{i=1}^{n_l} G(\pi_{\theta_i}) \hat{S}(\pi_{\theta_i}, \boldsymbol{v}_l, \rho)} + (1 - \alpha_l) \Big( \frac{\lambda_l}{n_l} \sum_{i=1}^{n_l} \Gamma(\theta_i) + (1 - \lambda_l) \hat{\eta}_l \Big), \tag{3}
$$

where  $v_l = m^{-1}(\hat{\eta}_l)$ .

We need to show the connection between [\(3\)](#page-2-0) and the ODE [\(8\)](#page-0-2) in Section [4.3.](#page-0-2) The first step is to rewrite [\(3\)](#page-2-0) to explicitly compare the sampling-based estimates to their true values. Or equivalently,

$$
\hat{\eta}_{l+1} - \hat{\eta}_{l} \n= \alpha_{l} \Big( \frac{\sum_{i=1}^{n_{l}} G(\pi_{\theta_{i}}) \hat{S}(\pi_{\theta_{i}}, \mathbf{v}_{l}, \rho) \Gamma(\theta_{i})}{\sum_{i=1}^{n_{l}} G(\pi_{\theta_{i}}) \hat{S}(\pi_{\theta_{i}}, \mathbf{v}_{l}, \rho)} - \hat{\eta}_{l} \Big) + (1 - \alpha_{l}) \lambda_{l} \Big( \frac{1}{n_{l}} \sum_{i=1}^{n_{l}} \Gamma(\theta_{i}) - \hat{\eta}_{l} \Big) \n= \alpha_{l} \Big( \frac{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta}) S(\pi_{\theta}, \mathbf{v}_{l}, \rho) \Gamma(\theta)]}{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta}) S(\pi_{\theta}, \mathbf{v}_{l}, \rho)]} - \hat{\eta}_{l} \Big) \n+ \alpha_{l} \Big( \frac{\sum_{i=1}^{n_{l}} G(\pi_{\theta_{i}}) \hat{S}(\pi_{\theta_{i}}, \mathbf{v}_{l}, \rho) \Gamma(\theta_{i})}{\sum_{i=1}^{n_{l}} G(\pi_{\theta_{i}}) \hat{S}(\pi_{\theta_{i}}, \mathbf{v}_{l}, \rho)} - \frac{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta}) S(\pi_{\theta}, \mathbf{v}_{l}, \rho) \Gamma(\theta)]}{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta}) S(\pi_{\theta}, \mathbf{v}_{l}, \rho)]} \Big) \n+ (1 - \alpha_{l}) \Big( \frac{\lambda_{l}}{n_{l}} \sum_{i=1}^{n_{l}} \Gamma(\theta_{i}) - \lambda_{l} \hat{\eta}_{l} \Big).
$$

<span id="page-2-2"></span>Define

$$
L_{l} = \frac{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_{l}, \rho)\Gamma(\theta)]}{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_{l}, \rho)]} - \hat{\eta}_{l}
$$
  
\n
$$
b_{l} = \frac{\sum_{i=1}^{n_{l}} G(\pi_{\theta_{i}})\hat{S}(\pi_{\theta_{i}}, \mathbf{v}_{l}, \rho)\Gamma(\theta_{i})}{\sum_{i=1}^{n_{l}} G(\pi_{\theta_{i}})\hat{S}(\pi_{\theta_{i}}, \mathbf{v}_{l}, \rho)} - \frac{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_{l}, \rho)\Gamma(\theta)]}{\mathbb{E}_{\mathbf{v}_{l}}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_{l}, \rho)]}
$$
(4)  
\n
$$
w_{l} = \frac{1 - \alpha_{l}}{\alpha_{l}} \Big( \frac{\lambda_{l}}{n_{l}} \sum_{i=1}^{n_{l}} \Gamma(\theta_{i}) - \lambda_{l} \hat{\eta}_{l} \Big),
$$

then [\(3\)](#page-2-0) can be rewritten as

<span id="page-2-1"></span>
$$
\hat{\eta}_{l+1} - \hat{\eta}_l = \alpha_l \Big( L_l + b_l + w_l \Big). \tag{5}
$$

Note that the first term in  $L_l$  coincides with  $\mathbb{E}_{v_l^*}[\Gamma(\theta)]$  where  $g_{v_l}^*$  is defined in Equation [\(3\)](#page-0-2) of Section [4.2.](#page-0-2) It holds by [\(6\)](#page-0-2) in Section [4.2](#page-0-2) that

$$
L_l = \mathbb{E}_{\mathbf{v}_l^*}[\Gamma(\theta)] - \hat{\eta}_l
$$
  
=  $\mathbb{E}_{\mathbf{v}_l^*}[\Gamma(\theta)] - m(\mathbf{v}_l)$   
=  $\frac{\partial}{\partial \mathbf{v}'} \log \mathbb{E}_{\mathbf{v}'}[G(\pi_\theta)S(\pi_\theta, \mathbf{v}_l, \rho)] \Big|_{\mathbf{v}' = \mathbf{v}_l}$ 

,

where the right hand side is the same as that of the ODE [\(7\)](#page-0-2) in Section [4.2.](#page-0-2)

We aim to show the connection between  $\{\eta_l\}_{l>0}$  (as defined in [\(3\)](#page-2-0) or [\(5\)](#page-2-1)) and the ODE [\(8\)](#page-0-2) using the following conclusion in stochastic approximation.

<span id="page-3-2"></span>**Theorem 1.** *(Theorem 1.2, [Benaim](#page-10-1) [\[1996\]](#page-10-1)) Let*  $H : \mathbb{R}^m \to \mathbb{R}^m$  *be a continuous vectorfield with unique integral curves. Let*  $\{w_n\}_{n\geq 0}$  *be the solution to*  $w_{n+1} - w_n = \gamma_n(H(w_n) + u_n + b_n)$ *, where* {γn}n≥<sup>0</sup> *is a decreasing gain sequence. Assume that*

- $\{\gamma_n\}_{n\geq 0}$  *is bounded.*
- $\lim_{n\to+\infty} b_n = 0$ .
- *For each*  $T > 0$ ,

$$
\lim_{n \to \infty} \Big( \sup_{k: 0 \le \tau_k - \tau_n \le T} \left| \left| \sum_{i=n}^{k-1} \gamma_i u_i \right| \right| \Big) = 0.
$$

*Then the limit set of*  $\{w_n\}_{n>0}$  *is a connected set internally chain-recurrent for the flow induced by* H*.*

We first show that  $\lim_{l\to\infty} b_l = 0$  where  $b_l$  is defined in [\(4\)](#page-2-2).

<span id="page-3-0"></span>**Lemma 3.** *Given Assumption* [\(2b\), \(2c\), \(2d\), \(2e\),](#page-0-2)  $\lim_{l\to\infty} b_l = 0$ , w.p.1.

In order to prove Lemma [3,](#page-3-0) we first show that the sample quantile is an unbiased estimate of the true quantile, which is stated in Lemma [4.](#page-3-1) Although we show the result for  $H$ , similar results apply for  $U$ .

<span id="page-3-1"></span>**Lemma 4.** Let  $\xi(\rho, \mathbf{v}_l)$  be the true  $(1 - \rho)$ -quantile of  $H(\pi_\theta)$  with  $\theta \sim f_{\mathbf{v}_l}$  and  $\hat{\xi}_l$  be a sample  $(1 - \rho)$ -quantile acquired from  $n_l$  i.i.d. samples. Given Assumption [\(2b\), \(2c\), \(2e\),](#page-0-2)  $\hat{\xi}_l - \xi(\rho, v_l) \to 0$  $as l \rightarrow \infty$  *w.p.l.* 

*Proof.* By Assumption [\(2e\),](#page-0-2)  $H(\pi) \in \mathcal{H} := [H_{min}, H_{max}]$  for all  $\pi \in \Pi$ . It can be verified that any true  $(1 - \rho)$ -quantile  $\xi(\rho, v_l)$  with  $\theta \sim f_{v_l}(\cdot)$  is an optimal solution of the following optimization problem [Homem-de Mello](#page-10-2) [\[2007\]](#page-10-2):

$$
\min_{\gamma \in \mathcal{H}} J_l(\gamma) := \mathbb{E}_{\mathbf{v}_l}[h(H(\pi_{\theta}), \gamma)]
$$
\n
$$
s.t. h(H(\pi_{\theta}), \gamma) = \begin{cases} (1 - \rho)(H(\pi_{\theta}) - \gamma), & \text{if } H(\pi_{\theta}) \ge \gamma, \\ \rho(\gamma - H(\pi_{\theta})), & \text{if } H(\pi_{\theta}) < \gamma. \end{cases}
$$

Similarly the sample  $(1 - \rho)$ -quantile  $\xi_l$  can be computed by minimizing

$$
\hat{J}_l(\gamma) := \frac{1}{n_l} \sum_{i=1}^{n_l} h(H(\pi_{\theta_i}), \gamma),
$$

where  $\{\theta_1, \ldots, \theta_{n_l}\}\$  are i.i.d. samples of distribution  $f_{\boldsymbol{v}_l}$ .

We first show that  $J_l(\gamma)$  uniformly converges to  $\hat{J}_l(\gamma)$  over H w.p.1, i.e.  $\sup_{\gamma \in H} |J_l(\gamma) - \hat{J}_l(\gamma)| \to 0$ as  $l \to \infty$  w.p.1.

Let  $\delta$  and  $r$  be two arbitrary scalars such that  $\delta > 0$  and  $r \le \frac{\delta}{3 \max(\rho, 1 - \rho)}$ . Let  $B(\gamma, r) := \{ \gamma' \in \mathcal{H} :$  $||\gamma - \gamma'|| \leq r$ } be the *r*-neighborhood of  $\gamma \in H$  within H. Since H is compact, there exists a finite

cover  $\mathcal{U} = \{h_1,\ldots,h_k\} \subset \mathcal{H}$  of  $\mathcal{H}$  such that  $\mathcal{H} \subseteq \bigcup_{i=1}^k B(h_i,r)$ . For each  $\gamma \in \mathcal{H}$ , let  $h(\gamma) \in \mathcal{U}$  be the closest component in H. By definition,  $\sup_{\gamma \in \mathcal{H}} ||\gamma - h(\gamma)|| \leq r$ . For any  $\gamma \in \mathcal{H}$ ,

$$
|J_l(\gamma) - J_l(h(\gamma))|
$$
  
\n
$$
= |\mathbb{E}_{\mathbf{v}_l}[h(H(\pi_{\theta}), \gamma)] - \mathbb{E}_{\mathbf{v}_l}[h(H(\pi_{\theta}), h(\gamma))]|
$$
  
\n
$$
\leq \max(\rho, 1 - \rho) ||\gamma - h(\gamma)|| \leq \frac{\delta}{3}.
$$
  
\n
$$
|\hat{J}_l(\gamma) - \hat{J}_l(h(\gamma))|
$$
  
\n
$$
= \frac{1}{n_l} |\sum_{i=1}^{n_l} \left( h(H(\pi_{\theta_i}), \gamma) - h(H(\pi_{\theta_i}), h(\gamma)) \right) |
$$
  
\n
$$
\leq \max(\rho, 1 - \rho) ||\gamma - h(\gamma)|| \leq \frac{\delta}{3}.
$$

As  $H(\cdot) \subseteq [H_{min}, H_{max}]$ , we can bound the probability that  $|J_l(h(\gamma)) - \hat{J}_l(h(\gamma))| > \delta$  for any  $\delta \geq 0$  by Hoeffding's inequality:

$$
Pr(|J_l(h(\gamma)) - \hat{J}_l(h(\gamma))| \ge \frac{\delta}{3})
$$
  

$$
\le 2 \exp(-\frac{2n_l \delta^2}{9|H_{max} - H_{min}|^2}).
$$

As card $(\mathcal{U}) = k < \infty$ , we can bound the probability that  $|J_l(h_i) - \hat{J}_l(h_i)| < \frac{\delta}{3}$  holds for all  $h_i \in \mathcal{U}$ with the union bound: δ

$$
Pr\left(\left(\max_{h_i \in \mathcal{U}} |J_l(h_i) - \hat{J}_l(h_i)|\right) \ge \frac{\delta}{3}\right)
$$
  

$$
\le \sum_{i=1}^k Pr\left(|J_l(h_i) - \hat{J}_l(h_i)| \ge \frac{\delta}{3}\right)
$$
  

$$
\le 2k \exp\left(-\frac{2n_l\delta^2}{9|H_{max} - H_{min}|^2}\right).
$$

Therefore with probability at least  $(1 - 2k \exp(-\frac{2n_l \delta^2}{9H} - H))$  $\frac{2n_l\delta^2}{9|H_{max}-H_{min}|^2})),$ 

$$
|J_l(\gamma) - \hat{J}_l(\gamma)| \le \frac{\delta}{3} + \frac{\delta}{3} + \frac{\delta}{3} = \delta
$$

holds uniformly for all  $\gamma \in \mathcal{H}$ . In other words,

$$
Pr(\sup_{\gamma \in \mathcal{H}} |J_l(\gamma) - \hat{J}_l(\gamma)| > \delta)
$$

$$
\leq 2k \exp(-\frac{2n_l \delta^2}{9|H_{max} - H_{min}|^2})).
$$

Therefore

$$
\sum_{l=1}^{\infty} Pr(\sup_{\gamma \in \mathcal{H}} |J_l(\gamma) - \hat{J}_l(\gamma)| > \delta)
$$
  

$$
\leq \sum_{l=1}^{\infty} 2k \exp(-\frac{2n_l \delta^2}{9|H_{max} - H_{min}|^2})) < \infty.
$$

By Assumption [\(2b\),](#page-0-2) the last inequality holds as  $n_l = \Theta(l^{\beta})$  with  $\beta > 0$ . By Borel-Cantelli Lemma,  $Pr(\sup_{\gamma \in \mathcal{H}} |J_l(\gamma) - \hat{J}_l(\gamma)| > \delta \text{ i.o.}) = 0$ . As the above proof holds for any  $\delta > 0$ ,  $\sup_{\gamma \in \mathcal{H}} |J_l(\gamma) - \hat{J}_l(\gamma)| \to 0$  as  $l \to \infty$  w.p.1. In other words,  $\hat{J}_l(\cdot)$  converges uniformly to  $J_l(\cdot)$  as  $l \rightarrow \infty$  w.p.1. Note that this uniform convergence holds whenever Assumption [\(2b\)](#page-0-2) and [\(2e\)](#page-0-2) hold.

Then we prove that  $\lim_{l \to +\infty} |\hat{\xi}_l - \xi(\rho, v_l)| = 0$ , w.p.1.

Since  $\sup_{\gamma \in \mathcal{H}} |J_l(\gamma) - \hat{J}_l(\gamma)| \to 0$  as  $l \to \infty$  w.p.1, there exists some  $L(\varepsilon) > 0$  for any  $\varepsilon > 0$  such that  $\sup_{\gamma \in \mathcal{H}} |J_l(\gamma) - \hat{J}_l(\gamma)| < \varepsilon$  holds for any  $l > L(\varepsilon)$ , w.p.1. Therefore with probability 1 and  $l > L(\varepsilon),$ 

$$
J_l(\hat{\xi}_l)-\varepsilon < \hat{J}_l(\hat{\xi}_l), \quad \hat{J}_l(\xi(\rho, \boldsymbol{v}_l)) < J_l(\xi(\rho, \boldsymbol{v}_l)) + \varepsilon.
$$

By definition of  $\xi(\rho, v_l)$  and  $\hat{\xi}_l$ , i.e.,  $\xi(\rho, v_l)$  minimizes  $J_l(\cdot)$  and  $\hat{\xi}_l$  minimizes  $\hat{J}_l(\cdot)$ , we get

$$
J_l(\xi(\rho, \boldsymbol{v}_l)) \leq J_l(\hat{\xi}_l), \quad \hat{J}_l(\hat{\xi}_l) \leq \hat{J}_l(\xi(\rho, \boldsymbol{v}_l)).
$$

Combining the above two equalities, we get

$$
J_l(\xi(\rho, \mathbf{v}_l)) - \varepsilon \leq J_l(\hat{\xi}_l) - \varepsilon < \hat{J}_l(\hat{\xi}_l)
$$
  

$$
\leq \hat{J}_l(\xi(\rho, \mathbf{v}_l)) < J_l(\xi(\rho, \mathbf{v}_l)) + \varepsilon.
$$

Therefore with probability 1,

$$
J_l(\xi(\rho, \boldsymbol{v}_l))-\varepsilon < J_l(\hat{\xi}_l) < J_l(\xi(\rho, \boldsymbol{v}_l)) + \varepsilon
$$

for sufficiently large *l*. In other words,  $J_l(\hat{\xi}_l) - J_l(\xi(\rho, \mathbf{v}_l)) \to 0$  as  $l \to +\infty$  w.p.1.

By Assumption [\(2c\),](#page-0-2) the  $(1 - \rho)$ -quantile of  $\{H(\pi_\theta) : \theta \sim f_\nu(\cdot)\}$  is unique for all  $\mathbf{v} \in \mathcal{V}$ . Therefore  $J_l(\gamma)$  is minimized with a unique  $\xi(\rho, v_l)$ . As  $J_l(\gamma)$  is also continuous in  $\gamma$ , for  $\varepsilon_l > 0$ that is small enough, there exists  $\delta_l(\varepsilon_l) > 0$  such that  $|J_l(\gamma) - J_l(\xi(\rho, v_l))| < \varepsilon_l$  if and only if  $||\xi(\rho, v_l) - \gamma|| < \overline{\delta}_l(\varepsilon_l)$ . Moreover,  $\overline{\delta}_l(\varepsilon_l) \to 0^+$  as  $\varepsilon_l \to 0^+$  for each l.

Assume that  $\hat{\xi}_l - \xi(\rho, v_l)$  does not converge to 0 w.p.1. Then  $\exists \bar{\delta} > 0$  such that  $Pr(\{|\hat{\xi}_l - \xi(\rho, v_l)| >$  $\overline{\delta}$  i.o.}) > 0. With positive probability, there exists a subsequence  $\{l_k\}_{k\geq 0} \in \mathbb{N}^{\infty}$  such that  $|\hat{\eta}_{k} - \xi(\rho, v_{l_k})| > \bar{\delta}$  for each  $k \in \mathbb{N}$  and  $\lim_{k \to \infty} J_{l_k}(\hat{\eta}_{l_k}) - J_{l_k}(\xi(\rho, v_{l_k})) = 0$ . We can further select a subsequence  $\{l_{k_j}\}_{j\geq 0}\subset \{l_k\}_{k\geq 0}$  such that  $|J_{l_{k_j}}(\hat{\eta}_{l_{k_j}})-J_{l_{k_j}}(\xi(\rho,\boldsymbol{v}_{l_{k_j}}))|<\frac{1}{2^j}$  and  $|\hat{\eta}_{l_{k_j}}-J_{l_{k_j}}(\eta)\rangle$  $|\xi(\rho, v_{l_{k_j}})| > \bar{\delta}$ . Since each  $\xi(\rho, v_{l_{k_j}})$  is unique, there exists  $j' \in \mathbb{N}$  such that  $\delta_{l_{k_{j'}}}(\frac{1}{2^{j_j}})$  $\frac{1}{2^{j'}}$ ) <  $\overline{\delta}$ , which contradicts our assumption that such a sequence  $\{l_k\}_{k\geq 0}$  exists.

 $\Box$ 

Therefore 
$$
\lim_{l \to +\infty} |\hat{\xi}_l - \xi(\rho, \boldsymbol{v}_l)| \to 0
$$
 w.p.1.

We can now give a proof to Lemma [3.](#page-3-0)

*Proof.* By Assumption [\(2e\),](#page-0-2)  $\inf_{\pi \in \Pi} G(\pi) > 0$ . By definition of  $(1 - \rho)$ -quantile, it holds for any  $v \in V$  that

$$
\mathbb{E}_{\boldsymbol{v}}[G(\pi_{\theta})S(\pi_{\theta},\boldsymbol{v},\rho)] \geq \inf_{\pi \in \Pi} G(\pi)\rho > 0.
$$

Similarly we can show

$$
\sum_{i=1}^{n_l} G(\pi_{\theta_i}) S(\pi_{\theta_i}, \boldsymbol{v}, \rho) \geq \inf_{\pi \in \Pi} G(\pi) > 0.
$$

There are two types of approximation involved in  $b_l$ : the first is to approximate  $\xi_H(\rho, \bm v_l)$  and  $\xi_U(\rho, \bm v_l)$ by  $\hat{\xi}_{H,l}$  and  $\hat{\xi}_{U,l}$ . The second is to approximate the expectation (e.g.  $\mathbb{E}_{v_l}[G(\pi_\theta)\hat{S}(\pi_{\theta_i}, v_l, \rho)\Gamma(\theta)]$ ) with sample mean (e.g.  $\frac{1}{n_l} \sum_{i=1}^{n_l} G(\pi_{\theta_i}) \hat{S}(\pi_{\theta_i}, \mathbf{v}_l, \rho) \Gamma(\theta_i)$ ).

As we have shown that  $\lim_{l\to\infty} |\xi_H(\rho, v_l) - \hat{\xi}_{H,l}| = 0$  w.p.1 and  $\lim_{l\to\infty} |\xi_U(\rho, v_l) - \hat{\xi}_{U,l}| = 0$ w.p.1 by Lemma [4,](#page-3-1) we can also get  $\lim_{l\to\infty} |S(\pi_\theta, \mathbf{v}_l, \rho) - \hat{S}(\pi_\theta, \mathbf{v}_l, \rho)| = 0$  w.p.1. We only need to consider the second part in this proof.

 $\Gamma(\cdot)$  is bounded as it is a continuous function defined over a compact set (by Assumption [\(2d\)\)](#page-0-2). By Assumption [\(2e\),](#page-0-2) both G and H are bounded over Π. By Remark [1,](#page-0-3)  $\delta(x \circ y)$  is bounded (by 1, to be specific) and Lipschitz continuous in both x and y. Let  $M > 0$  be a constant such that  $\sup_{\theta \in \Theta} |G(\pi_{\theta}) \Gamma(\theta)| \leq M$ . Therefore  $\lim_{l \to \infty}$  $\frac{1}{n_l}\sum_{i=1}^{n_l}G(\pi_{\theta_i})\hat{S}(\pi_{\theta_i},\boldsymbol{v}_l,\rho)\Gamma(\theta_i)\, \frac{1}{n_l}\sum_{i=1}^{n_l}G(\pi_{\theta_i})S(\pi_{\theta_i},\boldsymbol{v}_l,\rho)\Gamma(\theta_i)\Big| = 0$  w.p.1.

As  $G(\pi_{\theta})$ ,  $S(\pi_{\theta}, v_l, \rho)$ ,  $\Gamma(\theta)$  are all bounded for any  $\theta$  and  $\rho$ , there exist finite  $a, b$  such that  $a \leq G(\pi_{\theta})S(\pi_{\theta},\bm{v}_l,\rho)\Gamma(\theta) \leq b$  for any  $\theta \in \Theta$ . By Hoeffding's inequality, for any  $\varepsilon > 0$ 

$$
Pr\Big(\Big|\frac{1}{n_l}\sum_{i=1}^{n_l}G(\pi_{\theta_i})S(\pi_{\theta_i},\mathbf{v}_l,\rho)\Gamma(\theta_i)-\mathbb{E}_{\mathbf{v}_l}[G(\pi_{\theta})S(\pi_{\theta},\mathbf{v}_l,\rho)\Gamma(\theta)]\Big|\geq\varepsilon\Big)\n\leq 2\exp\Big(\frac{-2n_l\varepsilon^2}{(b-a)^2}\Big).
$$

By Assumption [\(2b\),](#page-0-2)  $n_l = \Theta(l^{\beta})$  and  $\beta > 0$ . Therefore for any  $\varepsilon > 0$ ,

$$
\sum_{l=1}^{\infty} Pr\left(\left|\frac{1}{n_l}\sum_{i=1}^{n_l} G(\pi_{\theta_i})S(\pi_{\theta_i},\mathbf{v}_l,\rho)\Gamma(\theta_i) - \mathbb{E}_{\mathbf{v}_l}[G(\pi_{\theta})S(\pi_{\theta},\mathbf{v}_l,\rho)\Gamma(\theta)]\right| \geq \varepsilon\right)
$$
  

$$
\leq \sum_{l=1}^{\infty} 2 \exp\left(\frac{-2n_l\varepsilon^2}{(b-a)^2}\right) < \infty.
$$

Then by Borel-Cantelli Lemma,

$$
\left| \frac{1}{n_l} \sum_{i=1}^{n_l} G(\pi_{\theta_i}) S(\pi_{\theta_i}, \mathbf{v}_l, \rho) \Gamma(\theta_i) - \mathbb{E}_{\mathbf{v}_l} [G(\pi_{\theta}) S(\pi_{\theta}, \mathbf{v}_l, \rho) \Gamma(\theta)] \right| \to 0, w.p.1.
$$

Therefore

$$
\frac{1}{n_l}\sum_{i=1}^{n_l}G(\pi_{\theta_i})S(\pi_{\theta_i},\boldsymbol{v}_l,\rho)\Gamma(\theta_i)-\mathbb{E}_{\boldsymbol{v}_l}[G(\pi_{\theta})S(\pi_{\theta},\boldsymbol{v}_l,\rho)\Gamma(\theta)]\to 0
$$

as  $l \rightarrow \infty$  w.p.1.

We can show that  $\frac{1}{n_l} \sum_{i=1}^{n_l} G(\pi_{\theta_i}) S(\pi_{\theta_i}, \mathbf{v}_l, \rho) - \mathbb{E}_{\mathbf{v}_l} [G(\pi_{\theta}) S(\pi_{\theta}, \mathbf{v}_l, \rho)] \to 0$  as  $l \to \infty$  w.p.1 in exactly the same way as above; it is a special case when  $\Gamma(\theta) = 1$  for all  $\theta \in \Theta$ .

By continuous mapping theorem, the facts that with probability 1,

$$
\mathbb{E}_{\mathbf{v}_l}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_l, \rho)] > 0, \ \forall \mathbf{v} \in \mathcal{V},
$$
\n
$$
\sum_{i=1}^{n_l} G(\pi_{\theta_i})\hat{S}(\pi_{\theta_i}, \mathbf{v}_l, \rho) > 0,
$$
\n
$$
\lim_{l \to \infty} \Big| \frac{1}{n_l} \sum_{i=1}^{n_l} G(\pi_{\theta_i})S(\pi_{\theta_i}, \mathbf{v}_l, \rho) \Gamma(\theta_i) - \mathbb{E}_{\mathbf{v}_l}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_l, \rho) \Gamma(\theta)] \Big| = 0,
$$
\n
$$
\lim_{l \to \infty} \frac{1}{n_l} \sum_{i=1}^{n_l} G(\pi_{\theta_i})\hat{S}(\pi_{\theta_i}, \mathbf{v}_l, \rho) - \mathbb{E}_{\mathbf{v}_l}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}_l, \rho)] = 0,
$$

guarantee that  $\lim_{l\to\infty} b_l = 0$ , w.p.1.

 $\Box$ 

### Now we restate Theorem [4.1](#page-0-2) and provide a proof.

**Theorem [4.1.](#page-0-2)** *If Assumptions*  $I - (2e)$  $I - (2e)$  *hold, the sequence*  $\{\hat{\eta}_l\}_{l\geq0}$  *in Step [1](#page-0-2)1 of Algorithm 1 converges to a connected internally chain recurrent set of* [\(8\)](#page-0-2) *as*  $l \rightarrow \infty$  *with probability 1.* 

*Proof.* We connect the sequence  $\{\hat{\eta}_l\}_{l\geq 0}$  to the ODE [\(8\)](#page-0-2) by applying Theorem [1.](#page-3-2) We need to verify that all sufficient conditions in [1](#page-3-2) hold properly. By [\(5\)](#page-2-1),  $\hat{\eta}_{l+1} - \hat{\eta}_l = \alpha_l (L_l + b_l + w_l)$ .

• By Assumption [\(2a\),](#page-0-2)  $\tilde{L}(v; \rho)$  is continuous in  $v \in int(V)$ . By Lemma [1,](#page-1-0)  $m^{-1}(\eta)$  is continuous in  $\eta$ . Therefore  $\tilde{L}(v;\rho)\Big|_{v=m^{-1}(\eta)}$  is continuous in  $\eta$ . [\(8\)](#page-0-2) has a unique integral curve by Assumption [\(2a\).](#page-0-2)

- By Assumption [\(2b\),](#page-0-2)  $\{\alpha_l\}_{l>0}$  is bounded and decreasing.
- By Lemma [3,](#page-3-0)  $\lim_{l\to\infty} b_l = 0$  w.p.1 with Assumption [\(2b\), \(2c\), \(2d\), \(2e\).](#page-0-2)
- Then we show that for any  $N \in \mathbb{N}^+$ ,  $\lim_{l \to \infty} \left( \sup_{k:n \leq k \leq n+N} \left| \left| \sum_{i=n}^k \alpha_i w_i \right| \right| \right) = 0$ .

Define  $M_n = \sum_{i=1}^n \alpha_i w_i$ . Then  $M_n = M_{n-1} + \alpha_n w_n$ . As the set  $\{\theta_i\}_{i=1}^{n_i}$  is generated i.i.d. with distribution  $f_{m^{-1}(\hat{\eta}_l)}(\cdot)$  and  $\hat{\eta}_l = \mathbb{E}_{m^{-1}(\hat{\eta}_l)}[\Gamma(\theta)],$ 

$$
\mathbb{E}[M_n | \sigma(M_1, ..., M_{n-1})]
$$
  
=  $M_{n-1} + \mathbb{E}_{n-1}(\hat{\eta}_l) [\frac{1}{n_l} \sum_{i=1}^{n_l} \Gamma(\theta_i) | M_{l-1}] - \hat{\eta}_l = M_{n-1}$ 

regardless of the value of  $\hat{\eta}_l$ . Therefore  $\{M_n\}_{n\geq 0}$  is a martingale. Note that  $w_i$  is independent on  $w_j$  if  $i \neq j$ , as all  $\theta$  are independently generated. Therefore  $\mathbb{E}[w_i^{\intercal} w_j] =$  $\mathbb{E}[w_i]$ <sup>T</sup> $\mathbb{E}[w_j] = 0$ .

$$
\mathbb{E}[||M_n||^2]
$$
  
\n
$$
= \mathbb{E}[M_n^{\mathsf{T}} M_n] = \mathbb{E}[(\sum_{i=1}^n \alpha_i w_i)^{\mathsf{T}} (\sum_{i=1}^n \alpha_i w_i)]
$$
  
\n
$$
= \sum_{i=1}^n \alpha_i^2 \mathbb{E}[w_i^{\mathsf{T}} w_i] + \sum_{i=1}^n \sum_{j \neq i} \alpha_i \alpha_j \mathbb{E}[w_i^{\mathsf{T}} w_j]
$$
  
\n
$$
= \sum_{i=1}^n \alpha_i^2 \mathbb{E}[w_i^{\mathsf{T}} w_i]
$$
  
\n
$$
= \sum_{i=1}^n \frac{(1 - \alpha_i)^2 \lambda_i^2}{n_i} \text{Cov}_{m^{-1}(\hat{\eta}_i)}[\Gamma(\theta)].
$$

As  $\Gamma(\theta)$  is continuous and the domain  $\Theta$  is compact, there exists  $0 < C < \infty$  such that  $Cov_{v}[\Gamma(\theta)] \leq C$  for any  $v \in V$ . Therefore by Assumption [\(2b\),](#page-0-2)

$$
\mathbb{E}[||M_n||^2] \le \sum_{i=1}^n C \frac{(1-\alpha_i)^2 \lambda_i^2}{n_i} = O(\sum_{l=1}^n \frac{1}{l^{\beta+2\lambda}}).
$$

By Assumption [\(2b\),](#page-0-2)  $\beta + 2\lambda > 1$  and thus  $\lim_{n\to\infty} \mathbb{E}[||M_n||^2] < \infty$ . As  $\{||M_n||^2\}$ increases monotonically, we know  $\sup_n \mathbb{E}[||M_n||^2] = \lim_{n \to \infty} \mathbb{E}[||M_n||^2] < \infty$ . Then by  $L_2$  martingale convergence theorem, there exists  $M_\infty$  such that  $M_n \to M_\infty$  w.p.1 and  $\mathbb{E}[||\tilde{M}_{\infty}||^2]<\infty.$ 

k

$$
\sup_{\{k:n\leq k\leq n+N\}} ||\sum_{i=n}^{k} \alpha_i w_i||
$$
  
= 
$$
\sup_{\{k:n\leq k\leq n+N\}} ||M_k - M_{n-1}|| \leq 2 \sup_{k\geq n} ||M_k||.
$$

Therefore

$$
0 \le \lim_{n \to \infty} \left( \sup_{\{k: n \le k \le n+N\}} \left| \left| \sum_{i=n}^{k} \alpha_i w_i \right| \right| \right)
$$
  

$$
\le \lim_{n \to \infty} \left( 2 \sup_{k \ge n-1} ||M_k|| \right) = 0
$$

for any finite  $N > 0$ .

Since all conditions in Theorem [1](#page-3-2) are satisfied, the limit set of sequence  $\{\hat{\eta}_l\}_{l\geq0}$  is a internally chain recurrent connected set for the flow induced by  $\overline{L}(\eta) := \frac{\mathbb{E}_{\mathbf{v}}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}, \rho)\Gamma(\theta)]}{\mathbb{E}_{\mathbf{v}}[G(\pi_{\theta})S(\pi_{\theta}, \mathbf{v}, \rho)]} \Big|_{\mathbf{v} = m^{-1}(\eta)} - \eta$  w.p.1. By [\(6\)](#page-0-2),  $\bar{L}(\eta) = \left(\frac{\partial}{\partial v} \log L(v;\rho)\right)^{\dagger} \Big|_{v=m^{-1}(\eta)}$ , which coincides with the right hand side of [\(8\)](#page-0-2).

#### Proof of Theorem [4.2](#page-0-2)

Now we restate Theorem [4.2](#page-0-2) and give a proof.

**Theorem [4.2.](#page-0-2)** Let  $\varphi : \mathcal{V} \to \mathbb{R}$  be any function such that  $\frac{\partial}{\partial v} \varphi(v) = \tilde{L}(v; \rho)$ . Any equilibrium  $\bar{v}^* \in int(V)$  of [\(9\)](#page-0-2) that is an isolated local maximum of  $\varphi(v)$  is locally asympototically stable.

*Proof.* The Lyapunov function we use is similar to that in [\[Joseph and Bhatnagar, 2016\]](#page-10-3):

$$
V(\boldsymbol{v}) := \varphi(\bar{\boldsymbol{v}}^*) - \varphi(\boldsymbol{v}),
$$

where  $\bar{v}^*$  is an isolated local maximum of  $\varphi(v)$  and v is in some neighborhood of  $\bar{v}^*$  such that  $\varphi(\bar{v}^*) \ge \varphi(v)$ , i.e.  $V(v) \ge 0$ . By previous analysis,  $\log \varphi(v)$  and  $V(v)$  are continuous in v. For the derivative:

$$
\frac{dV(\boldsymbol{v})}{dt} = -\frac{\partial \boldsymbol{v}}{\partial t} \frac{\partial \varphi(\boldsymbol{v})}{\partial \boldsymbol{v}} = -\Big(\tilde{L}(\boldsymbol{v};\rho)\Big)^{\sf T}(\mathrm{Cov}[\Gamma(\theta)])^{-1}\tilde{L}(\boldsymbol{v};\rho).
$$

As  $\text{Cov}_{\bm{v}}[\Gamma(\theta)]$  is positive definite for  $\bm{v} \in int(\mathcal{V})$ ,  $(\text{Cov}_{\bm{v}}[\Gamma(\theta)])^{-1}$  is also positive definite. Therefore  $\frac{\partial V(v)}{\partial t} \le 0$  in a neighborhood of  $v^*$  and  $\frac{\partial V(v)}{\partial t} = 0$  if and only if  $\tilde{L}(v; \rho) = 0$ , which guarantees that v is a stationary point of [\(9\)](#page-0-2). As  $\bar{v}^*$  is an isolated local maximum of  $\varphi(v)$ , it is the only stationary point in some neighborhood of  $\bar{v}^*$ . Therefore  $\frac{\partial V(v)}{\partial v} = 0$  if and only if  $v = \bar{v}^*$  (if v is in the neighborhood of  $v^*$ ) and  $\bar{v}^*$  is locally asymptotically stable.

In order to state the result we need to first introduce some definitions. By Assumption [\(2a\),](#page-0-2)  $Z :=$ If state the local we need to this introduce some definitions. By Assumption (Σα), B.<br>  $(\tilde{L}(v;\rho))^\mathsf{T}(\text{Cov}_v[\Gamma(\theta)])^{-1}$  is a continuous vector field defined on  $V \subset \mathbb{R}^{d_v}$  with unique integral curves. The *flow* of Z is the family of mappings  $\{\Phi_t(\cdot)\}_{t\in\mathbb{R}}$  defined on V by  $\frac{\partial \Phi_t(\mathbf{v})}{\partial t} = Z(\Phi_t(\mathbf{v}))$ such that  $\Phi_0(v) \equiv v$  and  $\Phi_t(\Phi_s(v)) \equiv \Phi_{t+s}(v)$  for any  $v \in V$ ,  $t, s \in \mathbb{R}$ .  $v \in V$  is an *equilibrium* if  $\Phi_t(\bm{v}) = \bm{v}$  for all t. A set  $\mathcal{V}' \subset \mathcal{V}$  is *positively invariant* under the flow  $\Phi$  if for all  $t \geq 0$ ,  $\Phi_t(\mathcal{V}')=\mathcal{V}'.$ 

#### Proof of Theorem [4.3](#page-0-2)

**Theorem [4.3.](#page-0-2)** *If all equilibria of* [\(9\)](#page-0-2) *are isolated, the sequence*  $\{v_l\}_{l>0}$  *derived by Algorithm [1](#page-0-2) converges toward an equilibrium of* [\(9\)](#page-0-2) *as*  $l \rightarrow \infty$  *with probability* 1.

*Proof.* Let  $\varphi$  be defined in the same way as in Theorem [4.2.](#page-0-2) We first show that  $\varphi$  is bounded over V. By definition of  $\tilde{L}(\mathbf{v};\rho)$  in [\(6\)](#page-0-2) of Section [4.2,](#page-0-2)  $\tilde{L}(\mathbf{v};\rho) = \frac{\mathbb{E}_{\mathbf{v}}[G(\pi_{\theta})S(\pi_{\theta},\mathbf{v},\rho)\Gamma(\theta)]}{L(\mathbf{v};\rho)} - m(\mathbf{v})$ . Since G has a positive lower bound (by Assumption [\(2e\)\)](#page-0-2) and  $\mathbb{E}_{v}[S(\pi_{\theta}, v', \rho)] \geq \rho$  for any  $v \in \mathcal{V}, L(v; \rho) \geq$  $\inf_{\pi \in \Pi} G(\pi)\rho > 0$ . Since  $\Gamma$  is continuous over  $\Theta$ ,  $\Theta$  and  $\mathcal V$  are compact (by Assumption [\(2d\)\)](#page-0-2),  $\Gamma(\theta)$  and  $m(\boldsymbol{v}) = \mathbb{E}_{\boldsymbol{v}}[\Gamma(\theta)]$  are both bounded. Since G is also bounded (by Assumption [\(2e\)\)](#page-0-2),  $\mathbb{E}_{\bm{v}}[G(\pi_{\theta})S(\pi_{\theta}, \bm{v}, \rho)\Gamma(\theta)]$  is also bounded over V for any  $\rho \in (0, 1)$ . Therefore  $\varphi$  is also bounded over  $\hat{V}$ .

Let  $\Phi$  be a flow induced by [\(9\)](#page-0-2) in Section [4.3](#page-0-2) and  $\Lambda$  be the set of all equilibria of (9). By definition, A is positively invariant under  $\Phi$ . Define  $V : \mathcal{V} \to \mathbb{R}^{\geq 0}$  as  $V(v) := \sup_{v' \in \mathcal{V}} \varphi(v') - \varphi(v)$ .  $\sup_{v \in V} \varphi(v') < \infty$  as  $\varphi$  is shown to be bounded in V. By definition of  $\Lambda$  and the proof of Theorem [4.2,](#page-0-2) the mapping  $t \mapsto V(\Phi_t(\mathbf{v}))$  is constant-valued for  $\mathbf{v} \in \Lambda$  and strictly decreasing for  $v \in int(V) \backslash \Lambda$ . Since we also assume that [\(9\)](#page-0-2) has only isolated equilibria and v is always in the interior of V (Assumption [\(2f\)\)](#page-0-2),  $\{v_l\}_{l>0}$  converges to an equilibrium of [\(9\)](#page-0-2) as  $l \to \infty$  with probability 1 by Corollary 3.3 in [\[Benaim, 1996\]](#page-10-1).  $\Box$ 

# Experiment Details

**Environment map and the local sensing model** The robot's state space is  $S = \{(x, y, \zeta) | x_{min} \leq$  $x \leq x_{max}, y_{min} \leq y \leq y_{max}, -\pi \leq \zeta < \pi$ , which contains the agent's position and orientation in the global coordinate; the input space is 2-dimensional:  $A = \{(v, \omega) | |v| \le v_{max}, |\omega| \le \omega_{max}\},$ in the grobal coordinate, the hip it space is 2-differential.  $A = \{(\nu, \omega)| |\nu| \le v_{max}, |\omega| \le \omega_{max}\}$ ,<br>which are linear and angular speed respectively. We assume that the robot can control v and  $\omega$  directly.

<span id="page-9-0"></span>

Figure 1: [\(1a\)](#page-9-0) Map of the car navigation example. There are one obstacle region (big grey rectangle), one goal region (small blue rectangle) and 10 randomly selected initial states (red circles pointing to the forward direction). Dotted lines are added to show  $x$  and  $y$  axes. [\(1b\)](#page-9-0) Illustrations of local features in the agent's local coordinate at one of the initial states, with  $n<sub>s</sub> = 5$ . Obstacle nodes, goal nodes and free nodes are labeled by black crosses, yellow plus signs and green triangles respectively. The goal direction (shown as the black arrow) is also included in local features.

There is a goal region G and a non-overlapping bad region B such that  $\mathcal{G}, \mathcal{B} \subset [x_{min}, x_{max}] \times$  $[y_{min}, y_{max}]$ . The map is shown in Figure [1a.](#page-9-0)

Since the robot has only local sensors, we use the following local sensing model instead of assuming the knowledge of the true state variables  $x, y$  and  $\zeta$ . For a given positive integer parameter  $n_s$ , we design a radial grid as  $n<sub>s</sub>$  circles in the agent's local coordinate. The difference between the diameters of adjacent circles is  $v_{max}\Delta t$ , where  $\Delta t$  is the sampling time. There are  $[2\pi/\omega_{max}]$  uniformly distributed nodes on each circle and the robot can measure the label for each node. A node is labeled 1 if it belongs to  $\mathcal{G}$ ; -1 if it belongs to  $\mathcal{B}$  and 0 otherwise. We also assume that the robot can sense the direction of the center of  $\mathcal G$  in its local coordinate without knowing the distance, so there are a total of  $(2 + n_s \lceil 2\pi/\omega_{max}\rceil)$  local features in total. The local features are illustrated in Figure [1b.](#page-9-0) In our experiment,  $\omega_{max} = \frac{\pi}{6}$ ,  $n_s = 5$ , so there are 62 local features as the inputs to the policy network. Note that the local sensor outputs are all discrete and only 2 features are continuous (the goal direction in the agent's local coordinate), so the problem is much simpler than a general continuous RL problem with the same number of continuous inputs.

**Algorithm parameters** In all experiments, we set  $F<sub>V</sub>$  as a class of multivariate Gaussian distributions with diagonal covariance matrices. The parameter space Θ contains all the parameters of the policy network. The policy space  $\Pi_{\Theta}$  is a set of deterministic stationary policies. Therefore the CCE trains a single neural network which takes states as inputs and output a single action. The two baseline algorithms TRPO [Schulman et al.](#page-10-4) [\[2015\]](#page-10-4) and CPO [Achiam et al.](#page-10-5) [\[2017\]](#page-10-5) take Gaussian policies, which takes states as inputs and outputs the mean and variance of the action distribution.

The policy networks for all experiments have two hidden layers of sizes (30, 30). The activation function for hidden layers is ReLU and that for the output layer is tanh. The length of sample trajectories are all 30. The same set of parameters are applied for all experiments. For CCE, we sample 40 different policies in each iteration. Each sampled policy is evaluated using 10 sample trajectories. The hyperparameter for selecting elite examples is  $\rho = 0.2$ . For both CPO and TRPO, the batch size is 6000, discount factor is 0.999, and the step size for trust region is 0.01. All the other parameters are used as default in the source code in rllab [Duan et al.](#page-10-6) [\[2016\]](#page-10-6).

The axes in the learning curve (Figure 1 in the paper) The x-axes in Figure 1 show the total number of sample trajectories for CCE or the total number of equivalent sample trajectories for TRPO and CPO. Assume that in each iteration of the CCE algorithm, we sample 40 policies (i.e.,  $n_l = 40$ ) and simulate 10 sample trajectories for each policy (Step 5 of Algorithm 1), then the total number of sample trajectories is 400 per iteration. If trajectory length is 30 and we sample 6000 new transitions in each iteration for CPO and TRPO, the number of equivalent sample trajectories is 200 per iteration. As we set the same trajectory length for all methods, the numbers of sample trajectories for all methods are comparable with each other.

The y-axes in Figure 1 show the *average* objective and constraint values of the learned policy. For CCE, the average values are computed with all rollout trajectories that are simulated with *all* the

policies sampled at the current iteration. Since we take the same number of rollout trajectories for each sample policy, the average value can be interpreted as the average performance of all sample policies at the current iteration. For CPO and TRPO, we simulate the current policy from exactly the same set of initial states and compute the average objective and constraint values for all trajectories. As a result, the comparison of different methods in Figure 1 is fair.

# **References**

- <span id="page-10-5"></span>J. Achiam, D. Held, A. Tamar, and P. Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, pages 22–31, 2017.
- <span id="page-10-1"></span>M. Benaim. A dynamical system approach to stochastic approximations. *SIAM Journal on Control and Optimization*, 34(2):437–472, 1996.
- <span id="page-10-6"></span>Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*, pages 1329– 1338, 2016.
- <span id="page-10-2"></span>T. Homem-de Mello. A study on the cross-entropy method for rare-event probability estimation. *INFORMS Journal on Computing*, 19(3):381–394, 2007.
- <span id="page-10-0"></span>J. Hu, P. Hu, and H. S. Chang. A stochastic approximation framework for a class of randomized optimization algorithms. *IEEE Transactions on Automatic Control*, 57(1):165–178, 2012.
- <span id="page-10-3"></span>A. G. Joseph and S. Bhatnagar. Revisiting the cross entropy method with applications in stochastic global optimization and reinforcement learning. In *ECAI 2016 - 22nd European Conference on Artificial Intelligence, 29 August-2 September 2016, The Hague, The Netherlands*, pages 1026–1034, 2016.
- <span id="page-10-4"></span>J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 1889–1897, 2015.