# HW 8

(one-pt grad estimator)

$q_t = \frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t))(r_t(x_t, a_t) - b_t(x_t))$     $b_t : X \to \mathbb{R}$; time varying baseline

$\nabla_a r_t(x_0, a_0) = \nabla_a r_t(x_0, a)\big|_{a=a_0}$

a) $r_t(x_t, \cdot)$; $\exists v_t(x_t) \in \mathbb{R}^d$ and $c_t(x_t) \in \mathbb{R}$ s.t.

$\nabla_a r_t(x_t, a) = c_t(x_t) + v_t(x_t)^T a$; $\underline{Prove}$ $\mathbb{E}_{a_t}(q_t) = v_t(x_t)$

$\mathbb{E}_{a_t}[q_t] = \mathbb{E}_{a_t}\left[\frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t))(r_t(x_t, a_t) - b_t(x_t))\right]$

$= \mathbb{E}_{a_t}\left[\frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t))(c_t(x_t) + v_t(x_t)^T a - b_t(x_t))\right] = \mathbb{E}_{a_t}\left[\frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t))(\underbrace{v_t(x_t)^T a + c_t(x_t) - b_t(x_t)}_{\to \text{ independent of action and } \mu_{z_t} = 0})\right]$

$= \mathbb{E}_{a_t}\left[\frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t)) v_t(x_t)^T a\right] = \frac{1}{\sigma^2}\mathbb{E}_{a_t}\left[(a_t - \mu_{\theta_t}(x_t)) v_t(x_t)^T a_t\right] = \frac{1}{\sigma^2}\mathbb{E}_{a_t}\left[z_t v_t(x_t)^T(z_t + \mu_{\theta_t}(x_t))\right] = \frac{1}{\sigma^2}\mathbb{E}_{a_t}\left[z_t v_t(x_t)^T z_t + z_t v_t(x_t)^T \mu_{\theta_t}(x_t)\right]$

$z_t = a_t - \mu_{\theta_t}(x_t)$

$= \frac{1}{\sigma^2}\left(\underbrace{\mathbb{E}_{a_t}[z_t v_t(x_t)^T z_t]}_{\text{independence}} + \mathbb{E}[z_t v_t(x_t)^T \mu_{\theta_t}(x_t)]\right) = \frac{1}{\sigma^2}\left(v_t(x_t)\mathbb{E}_{a_t}[z_t z_t^T] + v_t(x_t)\underbrace{\mathbb{E}_{a_t}[z_t]}_{\mu_{z_t}=0}\right) = \frac{1}{\sigma^2}v_t(x_t)\underbrace{\mathbb{E}_{a_t}[(a-\mu_\theta(x_t))(a-\mu_\theta(x_t))]}_{\sigma^2 \text{ by definition of covariance}}$

$= \frac{\sigma^2}{\sigma^2}v_t(x_t)^T = \boxed{v_t(x_t)^T}$

$\boxed{\therefore \mathbb{E}_{a_t}[q_t] = v_t(x_t)^T}$

$\underline{Prove:}$

b) $\|\mathbb{E}_{a_t}[q_t] - \nabla_a r_t(x_t, \mu_{\theta_t}(x_t))\| \leq \frac{d(d+2)(d+4)}{4} L\sigma$     $z_t = a_t - \mu_{\theta_t}(x_t)$     $z_t = a_t - \mu_{\theta_t}(x_t)$

Note: using Lemma 1: $\underbrace{\left|r_t(x_t, a) - [r_t(x_t, \mu_{\theta_t}(x_t)) + \nabla_a r_t(x_t, \mu_{\theta_t}(x_t))^T(a - \mu_{\theta_t}(x_t))]\right| \leq \frac{L}{2}\|a - \mu_{\theta_t}(x_t)\|^2}_{\to R(z)}$

for $r_t(x_t, \cdot)$ as provided in the question

$\underline{\mathbb{E}[q_t]}$

$q_t = \frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t))(r_t(x_t, a_t) - b_t(x_t))$

$\mathbb{E}[q_t] = \mathbb{E}\left[\frac{1}{\sigma^2}(a_t - \mu_{\theta_t}(x_t))(r_t(x_t, a_t) - b_t(x_t))\right] = \frac{1}{\sigma^2}\mathbb{E}\left[z_t(r_t(x_t, z_t + \mu_{\theta_t}) - b_t(x_t))\right] = \frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[z_t r_t(x_t, z_t + \mu_\theta) - z_t b_t(x_t)\right] = \frac{1}{\sigma^2}\left(\mathbb{E}_{z_t}\left[z_t r_t(x_t, z_t + \mu_\theta)\right] - \underbrace{\mathbb{E}_{z_t}[z_t b_t(x_t)]}_{\mu_{z_t}=0; \ b_t \text{ is a constant}}\right)$

$= \frac{1}{\sigma^2}\mathbb{E}\left[z_t r_t(x_t, z_t + \mu_{\theta_t})\right]$

$= \frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[z_t\underbrace{\left(r_t(x_t, \mu_{\theta_t}(x_t)) + \nabla_a r_t(x_t, \mu_{\theta_t}(x_t))^T(a_t - \mu_{\theta_t}(x_t)) + R(a_t)\right)}_{\text{first order Taylor expansion}}\right]$

$= \frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[z_t(r_t(x_t, \mu_{\theta_t}(x_t)) + \nabla_a r_t(x_t, \mu_{\theta_t}(x_t))^T z_t + z_t R(z_t)\right]$

$= \frac{1}{\sigma^2}\left(\underbrace{\mathbb{E}_{z_t}[z_t r_t(x_t, \mu_{\theta_t}(x_t))]}_{0} + \mathbb{E}_{z_t}\left[\underbrace{\nabla_a r_t(x_t, \mu_{\theta_t}(x_t))^T z_t}_{z_t}\right] + \mathbb{E}[z_t R(z_t)]\right)$

$$\frac{1}{\sigma^2}\left(\mathbb{E}_{z_t}\left[z_t \nabla_a r_t(x_t, \mu_\theta)^T z_t\right] + \mathbb{E}_{z_t}\left[z_t R(z_t)\right]\right) = \frac{1}{\sigma^2}\left(\nabla_a r_t(x_t, \mu_\theta)\mathbb{E}_{z_t}\left[z_t z_t^T\right] + \mathbb{E}_{z_t}\left[z_t R(z_t)\right]\right) = \frac{1}{\sigma^2}\left(\nabla_a r_t(x_t, \mu_\theta)\sigma^2 I + \mathbb{E}_{z_t}\left[z_t R(z_t)\right]\right)$$

$$= \nabla_a r_t(x_t, \mu_\theta) + \frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[z_t R(z_t)\right]$$

$$\mathbb{E}_{a_t}[g_t] = \nabla_a r_t(x_t, \mu_\theta) + \frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[z_t R(z_t)\right]$$

$$\left\|\mathbb{E}_{a_t}[g_t] - \nabla_a r_t(x_t, \mu_\theta(x_t))\right\| = \left\|\left(\frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[z_t R(z_t)\right]\right)\right\|$$

$$\mathbb{E}_{z_t}\left[\|z_t\| \|R(z_t)\|\right] \leq \frac{L}{2}\frac{1}{\sigma^2}\mathbb{E}_{z_t}\left[\|z_t\| \|z_t\|^2\right] = \frac{L}{\sigma^2}\mathbb{E}_{z_t}\left[\|z_t\|^3\right] \quad \text{using Lemma 1}$$

$$\text{let } w_t \sim N(0, I) \; ; \; z_t = \sigma w_t$$

$$\mathbb{E}_{z_t}\left[\|z_t\| \|R(z_t)\|\right] \leq \frac{L}{2\sigma^2}\mathbb{E}_{a_t}\left[\|\sigma w_t\|^3\right] = \frac{L}{2\sigma^2}\sigma^3\mathbb{E}_{w_t}\left[\|w_t\|^3\right] = \frac{L\sigma}{2}\sqrt{(\mathbb{E}_{w_t}[\|w_t\|^2])^3} = \frac{L\sigma}{2}\sqrt{d(d+2)(d+4)} \quad \text{using Lemma 2}$$

$$\boxed{\mathbb{E}_{z_t}\left[\|z_t\| \|R(z_t)\|\right] \leq L\sigma\sqrt{\frac{d(d+2)(d+4)}{4}}}$$

$$\boxed{\therefore \left\|\mathbb{E}_{a_t}[g_t] - \nabla_a r_t(x_t, \mu_{\theta_t}(x_t))\right\| \leq \sqrt{\frac{d(d+2)(d+4)}{4}}\, L\sigma}$$

$$\pi_\theta(a|x) = \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \exp\left(-\frac{\|a - \mu_\theta(x)\|^2}{2\sigma^2}\right)$$

c) $\underbrace{\frac{\pi_\theta(a_t|x_t)}{\pi_{\theta_t}(a_t|x_t)}\left(r_T(x_T,a) - b_T(x_t)\right) - \frac{1}{\eta} KL\left(\pi_\theta(\cdot|x_t), \pi_{\theta_t}(\cdot|x_t)\right)}_{} \cong \langle \mu_\theta(x_t) - \mu_{\theta_t}(x_t), g_t\rangle - \frac{1}{2\eta\sigma^2}\|\mu_\theta(x_t) - \mu_\theta(x_t)\|^2$

$\frac{\pi_\theta(a_t|x_t)}{\pi_{\theta_t}(a_t|x_t)}\left(r_T(x_t,a_t) - b_T(x_t)\right) - \frac{1}{\eta}\left(\log\frac{\sigma_{\theta_t}}{\sigma_{\pi_\theta}} + \frac{(\mu_{\pi_\theta} - \mu_{\theta_t})^2 + \sigma_\theta^2}{2\sigma_{\theta_t}} - \frac{1}{2}\right)$

$\frac{\frac{\exp\left(-\frac{\|a_t - \mu_\theta(x)\|^2}{2\sigma^2}\right)}{(2\pi\sigma_\theta^2)^{\frac{d}{2}}}}{\frac{\exp\left(-\frac{\|a_t - \mu_{\theta_t}(x)\|^2}{2\sigma^2}\right)}{(2\pi\sigma_\theta^2)^{\frac{d}{2}}}}\left(r_T(x_t,a_t) - b_T(x_t)\right) - \frac{1}{\eta}\left(\log\frac{\sigma_{\theta_t}}{\sigma_\theta} + \frac{(\mu_\theta - \mu_{\theta_t})^2 + \sigma_\theta^2}{2\sigma_{\theta_t}} - \frac{1}{2}\right)$

$\frac{\exp\left(-\frac{\|a_t - \mu_\theta(x_t)\|^2}{2\sigma_\theta^2}\right)(2\pi\sigma_\theta^2)^{\frac{d}{2}}}{\exp\left(-\frac{\|a_t - \mu_{\theta_t}(x_t)\|^2}{2\sigma_{\theta_t}^2}\right)(2\pi\sigma_\theta^2)^{\frac{d}{2}}}\left(r_T(x_t,a_t) - b_T(x_t)\right) - \frac{1}{\eta}\left(\log\frac{\sigma_{\theta_t}}{\sigma_\theta} + \frac{(\mu_\theta - \mu_{\theta_t})^2 + \sigma_\theta^2}{2\sigma_{\theta_t}} - \frac{1}{2}\right)$

$\exp\left(-\frac{\|a_t - \mu_\theta(x_t)\|^2}{2\sigma_\theta^2} + \frac{\|a_t - \mu_{\theta_t}(x_t)\|^2}{2\sigma_{\theta_t}^2}\right)\frac{\sigma_\theta^d}{\sigma_\theta^d}\left(r_T(x_t,a_t) - b_T(x_t)\right) - \frac{1}{\eta}\left(\log\frac{\sigma_{\theta_t}}{\sigma_\theta} + \frac{(\mu_\theta - \mu_{\theta_t})^2 + \sigma_\theta^2}{2\sigma_{\theta_t}} - \frac{1}{2}\right)$

$\exp\left(\frac{\|a_t - \mu_{\theta_t}(x_t)\|^2 - \|a_t - \mu_\theta(x_t)\|^2}{2\sigma_\theta^2}\right)\frac{\sigma_\theta^{\cancel{d}}{}^{1}}{\cancel{\sigma_\theta^d}}\left(r_T(x_t,a_t) - b_T(x_t)\right) - \frac{1}{\eta}\left(\cancel{\log\left(\frac{\sigma_{\theta_t}}{\sigma_\theta}\right)}^{0} + \frac{(\mu_\theta - \mu_{\theta_t})^2 + \sigma_\theta^2}{2\sigma_{\theta_t}} - \frac{1}{2}\right)$

$\rightarrow \|a_t - \mu_{\theta_t}(x_t)\|^2 = (a_t - \mu_{\theta_t}(x_t))^T(a_t - \mu_{\theta_t}(x_t)) = \|a_t\|^2 - 2a_t^T\mu_{\theta_t}(x_t) + \|\mu_{\theta_t}(x_t)\|^2$

<u>focused on first term</u>

$\exp\left(\frac{\overbrace{\|a_t\|^2 - 2a_t^T\mu_{\theta_t}(x_t) + \|\mu_{\theta_t}(x_t)\|^2 - \|a_t\|^2 + 2a_t^T\mu_\theta(x_t) - \|\mu_\theta(x_t)\|^2}}{2\sigma_\theta^2}\right)\left(r_T(x_t,a_t) - b_T(x_t)\right) - \ldots$

$\exp\left(\frac{1}{2\sigma_\theta^2}\left(2a_t^T(\mu_\theta(x_t) - \mu_{\theta_t}(x_t)) + \|\mu_{\theta_t}(x_t)\|^2 - \|\mu_\theta(x_t)\|^2\right)\right)\left(r_T(x_t,a_t) - b_T(x_t)\right) - \ldots$

let $\Delta\mu = \mu_\theta(x_t) - \mu_{\theta_t}(x_t)$

$\rightarrow \|\mu_\theta(x_t)\|^2 = \|\Delta\mu + \mu_{\theta_t}(x_t)\| = (\Delta\mu + \mu_{\theta_t}(x_t))^T(\Delta\mu + \mu_{\theta_t}(x_t)) = \|\Delta\mu\|^2 + 2\Delta\mu^T\mu_{\theta_t}(x_t) + \|\mu_{\theta_t}(x_t)\|^2$

$\exp\left(\frac{1}{2\sigma_\theta^2}\left(2a_t^T\Delta\mu + \cancel{\|\mu_{\theta_t}(x_t)\|^2} - \|\Delta\mu\|^2 - 2\Delta\mu^T\mu_{\theta_t}(x_t) - \cancel{\|\mu_{\theta_t}(x_t)\|^2}\right)\right)\left(r_T(x_t,a_t) - b_T(x_t)\right) - \ldots$

$\exp\left(\frac{1}{2\sigma_\theta^2}\left(2(a_t^T - \mu_{\theta_t}(x_t))\Delta\mu - \|\Delta\mu\|^2\right)\right)\left(r_T(x_t,a_t) - b_T(x_t)\right) - \ldots$

$$\left(r_f(x_t, a_t) - b_f(x_t)\right) + \frac{1}{2\sigma_\theta^2}\left(2(a_t - \mu_{\theta_f}(x_t))^T \Delta\mu - \|\Delta\mu\|^2\right)\left(r_f(x_t, a_t) - b_f(x_t)\right) - \cdots$$

$$\left(r_f(x_t, a_t) - b_f(x_t)\right) + \frac{2(a_t - \mu_{\theta_f}(x_t))^T \Delta\mu \left(r_f(x_t, a_t) - b_f(x_t)\right)}{2\sigma_\theta^2} - \frac{\left(r_f(x_t, a_t) - b_f(x_t)\right)\|\Delta\mu\|^2}{2\sigma_\theta^2} - \cdots$$

<span style="color:red">Constant w.r.t. θ so it doesn't affect argmax</span>

$$\underbrace{\left(r_f(x_t, a_t) - b_f(x_t)\right)}_{} + \Delta\mu\underbrace{\left(\frac{1}{\sigma_\theta^2}(a_t - \mu_{\theta_f}(x_t))^T\left(r_f(x_t, a_t) - b_f(x_t)\right)\right)}_{g_t} - \underbrace{\frac{1}{2\sigma_\theta^2}\left(r_f(x_t, a_t) - b_f(x_t)\|\Delta\mu\|^2\right)}_{} - \cdots$$

<span style="color:red">Since Δμ is small, it grows smaller than the linear term</span>

$$\langle \Delta\mu, g_t\rangle - \cdots = \underbrace{\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle}_{\text{first term}} - \underbrace{\cdots}_{KL}$$

## KL Portion

$$\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle - \frac{1}{\eta}\left(\frac{(\mu_\theta - \mu_{\theta_f})^2 + \sigma_\theta^2}{2\sigma_{\theta_f}^2} - \frac{1}{2}\right)$$

<span style="color:red">since $\theta_t \approx \theta_{t+1}$</span>

$$\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle - \frac{1}{\eta}\left(\frac{(\mu_\theta - \mu_{\theta_f})^2}{2\sigma_{\theta_f}^2} + \frac{\sigma_\theta^2}{2\sigma_{\theta_f}^2} - \frac{1}{2}\right)$$

$$\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle - \frac{1}{\eta}\left(\frac{(\mu_\theta - \mu_{\theta_f})^2}{2\sigma_{\theta_f}^2} + \frac{1}{2}\overset{0}{\cancel{- \frac{1}{2}}}\right)$$

<span style="color:red">same as $(\mu_\theta - \mu_{\theta_f})^2$ ; just a vectorized version</span>

$$\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle - \frac{1}{\eta}\left(\frac{(\mu_\theta - \mu_{\theta_f})^2}{2\sigma_{\theta_f}^2}\right) = \langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle - \frac{1}{2\eta\sigma_{\theta_f}^2}\|\mu_\theta - \mu_{\theta_f}\|^2$$

<span style="color:red">some constant unrelated to θ</span>

$$\therefore \arg\max_\theta \left(\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t)\rangle - \frac{1}{2\eta\sigma_{\theta_f}^2}\|\mu_\theta - \mu_{\theta_f}\|^2 + C\right) \cong \arg\max_\theta \left(\langle \mu_\theta(x_t) - \mu_{\theta_f}(x_t), g_t\rangle - \frac{1}{2\eta\sigma_{\theta_f}^2}\|\mu_\theta - \mu_{\theta_f}\|^2\right)$$

d) **PG**

$$\theta_{t+1} \leftarrow \theta_t + \eta \nabla_\theta \log \pi_\theta(a_t|x_t)\big|_{\theta=\theta_t} (r_t(x_t,a_t) - b_t(x_t))$$

$$\theta_{t+1} \leftarrow \text{argmax}_\theta \left\{ \langle \mu_\theta(x_t) - \mu_{\theta_t}(x_t), g_t \rangle - \frac{1}{2\eta}\|\theta - \theta_t\|^2 \right\}$$

$$\theta_t + \eta \nabla_\theta \log \pi_\theta(a_t|x_t)\big|_{\theta=\theta_t} (r_t(x_t,a_t) - b_t(x_t)) \cong \text{argmax}_\theta \left\{ \langle \mu_\theta(x_t) - \mu_{\theta_t}(x_t), g_t \rangle - \frac{1}{2\eta}\|\theta - \theta_t\|^2 \right\}$$

$$\theta_t + \eta \nabla_\theta \log\left( \frac{\exp\left(-\frac{\|a_t - \mu_\theta(x)\|^2}{2\sigma^2}\right)}{(2\pi\sigma^2)^{\frac{d}{2}}} \right)\bigg|_{\theta=\theta_t} (r_t(x_t,a_t) - b_t(x_t))$$

$$\nabla_\theta \log\left( \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \exp\left(-\frac{\|a_t - \mu_\theta(x_t)\|^2}{2\sigma^2}\right)\right) = \nabla_\theta\left( \ln\left(\frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}}\right) + \ln\left(\exp\left(-\frac{\|a_t - \mu_\theta(x_t)\|^2}{2\sigma^2}\right)\right)\right) = \nabla_\theta\left( \underbrace{-\frac{d}{2}\ln(2\pi\sigma^2)}_{\text{constant}} + -\frac{\|a_t - \mu_\theta(x_t)\|^2}{2\sigma^2}\right)$$

$$= \nabla_\theta\left(-\frac{\|a_t - \mu_\theta(x_t)\|^2}{2\sigma^2}\right) = \nabla_\theta\left(-\frac{1}{2\sigma^2}(a_t - \mu_\theta(x_t))^T(a_t - \mu_\theta(x_t))\right) = -\frac{1}{2\sigma^2}\nabla_\theta\left(\overset{0}{\|a_t\|^2} - 2a_t^T\mu_\theta(x_t) + \|\mu_\theta(x_t)\|^2\right)$$

$$= -\frac{1}{2\sigma^2}\left(-2a_t^T\frac{\partial \mu_\theta(x_t)}{\partial\theta} + 2\mu_\theta(x_t)\frac{\partial \mu_\theta}{\partial\theta}\right) = \frac{1}{\sigma^2}\left(a_t^T\frac{\partial \mu_\theta(x_t)}{\partial\theta} - \mu_\theta(x_t)\frac{\partial \mu_\theta(x_t)}{\partial\theta}\right)\bigg|_{\theta=\theta_t} = \frac{1}{\sigma^2}(a_t - \mu_\theta(x_t))\underbrace{\frac{\partial \mu_\theta(x_t)}{\partial\theta}}_{\hookrightarrow\ J_{\theta_t}\ \text{for Jacobian}}\bigg|_{\theta=\theta_t}$$

$$\theta_t + \frac{\eta}{\sigma^2}(a_t - \mu_\theta(x_t))\, J_{\theta_t}\,(r_t(x_t,a_t) - b_t) = \theta_t + \eta \underbrace{\frac{1}{\sigma^2}(a_t - \mu_\theta(x_t))(r_t(x_t,a_t) - b_t)}_{g_t}\, J_{\theta_t}$$

$$\underline{\theta_{t+1} = \theta_t + \eta g_t J_{\theta_t}} \quad \text{PG update}$$

· · ‒ ‒ ‒ ‒‒ ·· **Right Side of Term**

$$\nabla_\theta\left(\langle \mu_\theta(x_t) - \mu_{\theta_t}(x_t), g_t \rangle - \frac{1}{2\eta}\|\theta - \theta_t\|^2\right) = 0 \qquad \text{Note: } \theta = \theta_{t+1}$$

$$\nabla_\theta\left((\mu_\theta(x_t) - \mu_{\theta_t}(x_t))^T g_t - \frac{1}{2\eta}(\theta - \theta_t)^T(\theta - \theta_t)\right) = 0$$

$$\nabla_\theta\left((\theta - \theta_t)^T \nabla_\theta \mu_{\theta_t}(x_t) g_t - \frac{1}{2\eta}(\theta - \theta_t)^T(\theta - \theta_t)\right) = 0 \Rightarrow \nabla_\theta\left((\theta - \theta_t)^T\frac{\partial \mu_\theta(x_t)}{\partial\theta} g_t - \frac{1}{2\eta}(\theta - \theta_t)^T(\theta - \theta_t)\right) = 0 \Rightarrow \nabla_\theta\left((\theta - \theta_t)^T J_{\theta_t} g_t - \frac{1}{2\eta}(\theta - \theta_t)^T(\theta - \theta_t)\right)$$

$$= (1 - 0)^T J_{\theta_t} g_t - \frac{1}{2\eta}\left(\overset{0}{\|\theta\|^2} - 2\theta^T\theta_t + \|\theta_t\|^2\right)$$

$$= J_{\theta_t} g_t - \frac{1}{2\eta}(2\theta - 2\theta_t) = J_{\theta_t} g_t - \frac{1}{\eta}(\theta - \theta_t) = 0$$

$$J_{\theta_t} g_t - \frac{1}{\eta}(\theta_{t+1} - \theta_t) = 0 \longrightarrow J_{\theta_t} g_t = \frac{1}{\eta}(\theta_{t+1} - \theta_t) \longrightarrow \eta J_{\theta_t} g_t = \theta_{t+1} - \theta_t \longrightarrow \boxed{\theta_{t+1} = \theta_t + \eta J_{\theta_t} g_t}$$

The PG update is equivalent to its argmax version since $\theta_t + \eta \underbrace{J_{\theta_t} g_t}_{\text{PG Update}} = \theta_t + \eta \underbrace{J_{\theta_t} g_t}_{\text{equivalence}}$