

$$\boxed{1.1} \quad \hat{q} = \underset{q}{\operatorname{argmin}} D_{\text{KL}}(q(z) \parallel p_{\theta}(z|n)) = \underset{q}{\operatorname{argmin}} \mathbb{E} \left[ \lg \frac{q(z)}{p_{\theta}(z|n)} \right]$$

$$= \underset{q}{\operatorname{argmin}} \mathbb{E} [\lg q(z)] - \mathbb{E} [\lg (p_{\theta}(n|z) p_{\theta}(z))] + \underbrace{\mathbb{E} [\lg p_{\theta}(n)]}_{\lg p_{\theta}(n) \rightarrow \text{const}}$$

$$\begin{aligned} \rightarrow \hat{q} &= \underset{q}{\operatorname{argmax}} \mathbb{E} [\lg (p_{\theta}(n|z) p_{\theta}(z)) - \lg q(z)] \\ &= \underset{q}{\operatorname{argmax}} \mathbb{E} [\lg p_{\theta}(n|z)] - D_{\text{KL}}(q(z) \parallel p_{\theta}(z)) = \underset{q}{\operatorname{argmax}} \text{ELBO} \end{aligned}$$

$$\lg p_{\theta}(n) = \lg \int p_{\theta}(n, z) dz = \lg \int p_{\theta}(n, z) \frac{q(z)}{q(z)} dz$$

$$\underbrace{\mathbb{E}_{q(z)} \left[ \frac{p(n, z)}{q(z)} \right]}$$

$$\xrightarrow{\text{Jensen's inequality}} \lg p_{\theta}(n) \geq \mathbb{E}_{q(z)} \left[ \lg \frac{p(n, z)}{q(z)} \right] = \mathbb{E}_{q(z)} \left[ \lg \frac{p(n|z) p(z)}{q(z)} \right] = \underbrace{\mathbb{E}_{q(z)} [\lg p_{\theta}(n|z)] - D_{\text{KL}}(q(z) \parallel p_{\theta}(z))}_{\text{ELBO}}$$

$\boxed{1.2}$  ابتدا لاس را به صورت مقابل در نظر میگیریم و توسط الگوریتم زیر آنرا کمینه میگیریم

$$L(x_i; \psi, \theta) = -\text{ELBO}(x_i; q, \theta)$$

$\theta, \psi \leftarrow \text{random init}$

while  $L(n; \psi, \theta)$  has not converged :

for  $i \in 1 \dots N$ :

$$\psi_i \leftarrow \underset{\psi_i}{\operatorname{argmax}} L(\psi_i, \theta; x_i)$$

calc  $L(n; \psi, \theta)$

$$\theta \leftarrow \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N L(x_i; \psi_i, \theta)$$

return  $\psi, \theta$

1.3.a

stochastic VI

$\theta, \psi \leftarrow$  random init

while  $L(x; \psi, \theta)$  has not converged:

for  $B \in D$ :

for  $x_j \in B$ :

$l_j \leftarrow \text{calc } L(x_j, \psi_j, \theta)$

$g_j \leftarrow \text{calc } \nabla_{\psi_j} L(x_j, \psi_j, \theta)$

$h_j \leftarrow \text{calc } \nabla_{\theta} L(x_j, \psi_j, \theta)$

$$\hat{g} = \frac{1}{M} \sum_{j=1}^M g_j, \quad \hat{h} = \frac{1}{M} \sum_{j=1}^M h_j$$

for  $x_j \in B$ :

$$\bar{g} = [I(\psi_j)]^{-1} \hat{g}$$

$\psi_j \leftarrow$  update using  $\bar{g}$

$$\bar{h} = [I(\theta)]^{-1} \hat{h}$$

$\theta \leftarrow$  update using  $\bar{h}$

$$\hat{L}(B; \psi, \theta) = \sum_{j=1}^M l_j$$

$$L(x; \psi, \theta) \simeq \frac{N}{M} \hat{L}(B, \psi, \theta)$$

return  $\psi, \theta$

:  $p_i$  و  $m_i$  من  $B$  و  $m_i$  batch size = M و data loader  $D$  و  $D$  من

$$I(\psi_i) \triangleq \mathbb{E}_g \left[ \nabla_{\psi_i} \log q(z|x_i; \psi_i) \nabla_{\psi_i} \log q(z|x_i; \psi_i)^T \right]$$

$$I(\theta) \triangleq \mathbb{E}_{p(x_i, z; \theta)} \left[ \nabla_{\theta} \log p(x_i, z; \theta) \nabla_{\theta} \log p(x_i, z; \theta)^T \right]$$

عند  $\psi_{1:N}$  و  $\psi$  و  $\theta$  و  $\theta$

1.3.6

(gradient ascent) Amortized VI

$\theta, \phi \leftarrow \text{random init}$

while  $L(\pi; \phi, \theta)$  has not converged:

for  $B \in \mathcal{D}$ :

for  $\pi_j \in B$ :

$\ell_j \leftarrow \text{calc } L(\pi_j; \phi, \theta)$

$g_j \leftarrow \text{calc } \nabla_{\phi} L(\pi_j; \phi, \theta)$

$h_j \leftarrow \text{calc } \nabla_{\theta} L(\pi_j; \phi, \theta)$

$$\hat{g} = \frac{1}{M} \sum_{j=1}^M g_j, \quad \hat{h} = \frac{1}{M} \sum_{j=1}^M h_j$$

$\phi \leftarrow \text{update using } \hat{g}$

$\theta \leftarrow \text{update using } \hat{h}$

$$\hat{L}(B; \phi, \theta) = \sum_{j=1}^M \ell_j$$

$$L(\pi; \phi, \theta) \simeq \frac{N}{M} \hat{L}(B; \phi, \theta)$$

return  $\phi, \theta$

$$1.4 \quad ELBO = \mathbb{E}_{q_\phi(z)} \left[ \lg \frac{1}{(2\pi)^{\frac{k}{2}} \sigma^k} - \frac{1}{2\sigma^2} \|x - f_\theta(z)\|^2 \right] \\ - \frac{1}{2} \left[ k\sigma_\phi^2(x) + \|\mu_\phi(x)\|^2 - k + \lg \frac{1}{\sigma_\phi^{2k}(x)} \right] \\ - 2k \lg \sigma_\phi(x)$$

$$\text{loss} = -ELBO$$

$$\rightarrow \nabla_\theta \text{loss} = \nabla_\theta \mathbb{E}_{q_\phi(z)} \left[ -\frac{1}{2\sigma^2} \|x - f_\theta(z)\|^2 \right] = -\frac{1}{2\sigma^2} \mathbb{E} \left[ \underbrace{\nabla_\theta \|x - f_\theta(z)\|^2}_{2(f'_\theta(z) - x) f'_\theta(z)} \right]$$

برای  $\nabla_\phi$  اما نمی‌توانیم گزاردن را داخل  $\mathbb{E}$  ببرد. بنابراین از reparameterization استفاده می‌کنیم:

$$\nabla_\phi \text{loss} = \nabla_\phi \mathbb{E}_{q_\phi(z)} \left[ -\frac{1}{2\sigma^2} \|x - f_\theta(z)\|^2 \right] - \frac{k}{2} \nabla_\phi \sigma_\phi^2(x) + \nabla_\phi \|\mu_\phi(x)\|^2 - 2k \nabla_\phi \lg \sigma_\phi(x)$$

$$\frac{\hat{z} = \mu + \sigma \varepsilon}{\varepsilon \sim q(\varepsilon)} \quad \frac{-1}{2\sigma^2} \mathbb{E}_{q(\varepsilon)} \left[ \underbrace{\nabla_\phi \|x - f_\theta(\hat{z})\|^2}_{2(f'_\theta(\hat{z}) - x) f'_\theta(\hat{z}) \hat{z}'} \right] - k \sigma_\phi(x) \sigma'_\phi(x) + 2\mu_\phi(x) \mu'_\phi(x) - 2k \frac{\sigma'_\phi(x)}{\sigma_\phi(x)}$$

$\varepsilon$  مستقل از  $\phi$  است

با توجه به مشتق‌های به دست آمده می‌توان این مقادیر را در الگوریتم‌های قبلی جایگذاری کرد یا اگر نتایج را داشته باشیم شاید بتوان مشتق‌ها را برابر 0 گذاشت و باینک گام به جواب رسید

**2.a**  $q(z_t | x) = \mathcal{N}(a_t x, \sigma_t^2 I) \rightarrow z_t = a_t x + \sigma_t \varepsilon$

$\varepsilon \sim \mathcal{N}(0, I)$

$\rightarrow \text{var}[z_t] = \underbrace{\text{var}[a_t x]}_{a_t^2 I} + \underbrace{\text{var}[\sigma_t \varepsilon]}_{\sigma_t^2 I} = (1 - \sigma_t^2) I + \sigma_t^2 I = I$

بنابراین اگر  $a_t^2 = 1 - \sigma_t^2$  داریم ثابت می ماند

**2.b**

$q(z_s, z_t) = q(z_t | z_s) q(z_s) \rightarrow q(z_s, z_t | x) = \underbrace{q(z_t | z_s, x)}_{\text{چون } x \text{ معلوم است! داریم } z_s}$

$\rightarrow q(z_s, z_t | x) = q(z_t | z_s) q(z_s | x)$

$z_t$  از  $x$  مستقل است

**2.c**

**2.a** طبق

$z_t = a_t x + \sigma_t \varepsilon$

$z_s = a_s x + \sigma_s \varepsilon$

$x = \frac{z_s - \sigma_s \varepsilon}{a_s}$

$\rightarrow z_t | z_s = a_t \frac{z_s - \sigma_s \varepsilon}{a_s} + \sigma_t \varepsilon$

$E[z_t | z_s] = \frac{a_t}{a_s} z_s$

$\text{var}[z_t | z_s] = \left(\frac{a_t}{a_s}\right)^2 \underbrace{\text{var}[z_s]}_0 - \left(\frac{a_t}{a_s}\right)^2 \sigma_s^2 I + \sigma_t^2 I$   
 $= \left(\sigma_t^2 - a_{t|s}^2 \sigma_s^2\right) I$

چون  $z_s$  را به صورت given داریم پس داریم آن صفر است (deterministically آنرا می دانیم)

بنابراین  $q(z_t | z_s)$  یک گاوسی است که میانگین و واریانس آنرا بدست آوردیم

$q(z_t | z_s) = \mathcal{N}(a_{t|s} z_s, \sigma_{t|s}^2 I)$

2.d  $q(z_s | z_t, x) = \frac{q(z_t | z_s, x) q(z_s | x)}{q(z_t | x)} \propto q(z_t | z_s) q(z_s | x)$

این دو گامی را از بخش های قبل داریم. می دانیم که برای ضرب دو گامی فنول زیر بردار است:

$$q(z_s | x) \times q(z_t | z_s) = \mathcal{N}(\tilde{\mu}, \tilde{\Sigma}^2 I)$$

$$\tilde{\Sigma}^{-2} = \Sigma_A^{-2} + a^2 \Sigma_B^{-2}$$

$$\tilde{\mu} = \tilde{\Sigma}^2 (\Sigma_A^{-2} \mu_A + a \Sigma_B^{-2} z_t)$$

بنابراین برای مسئله مان داریم:

$$\tilde{\Sigma}^{-2} = \Sigma_s^{-2} + a_{t|s}^2 \Sigma_{t|s}^{-2} = \frac{1}{\Sigma_s^2} + \frac{a_{t|s}^2}{\Sigma_{t|s}^2} \rightarrow \tilde{\Sigma}^2 = \frac{\Sigma_{t|s}^2 \Sigma_s^2}{\Sigma_{t|s}^2 + \Sigma_s^2 a_{t|s}^2} = \frac{\Sigma_{t|s}^2 \Sigma_s^2}{\Sigma_t^2}$$

$$\tilde{\mu} = \frac{\Sigma_{t|s}^2 \Sigma_s^2}{\Sigma_t^2} \left( \frac{a_s x}{\Sigma_s^2} + \frac{a_{t|s} z_t}{\Sigma_{t|s}^2} \right) = \frac{a_{t|s} \Sigma_s^2}{\Sigma_t^2} z_t + \frac{a_s \Sigma_{t|s}^2}{\Sigma_t^2} x$$

2.e  $\Sigma_0 = \Sigma_1$  است پس  $\Sigma_Q^2 I$  برابر  $p_\theta(z_s | z_t)$  داریم سوال دارا این  $p_\theta(z_s | z_t)$  است

$$D_{KL}(q(z_s | z_t, x) \| p_\theta(z_s | z_t)) = \frac{1}{2} \left( \text{tr} \left[ \underbrace{(\Sigma_Q^2 I)^{-1} \Sigma_Q^2 I}_I \right] + (\mu_\theta - \mu_Q)^T (\Sigma_Q^2 I)^{-1} (\mu_\theta - \mu_Q) - d + \underbrace{\log \frac{\det(\Sigma_Q^2 I)}{\det(\Sigma_Q^2 I)}}_0 \right)$$

$$= \frac{1}{2 \Sigma_Q^2} \left( (\mu_\theta - \mu_Q)^T (\mu_\theta - \mu_Q) \right) = \frac{1}{2 \Sigma_Q^2} \|\mu_Q - \mu_\theta\|_2^2$$

$$\begin{aligned}
 \boxed{2.f} \quad D_{KL} &= \frac{1}{2\sigma_Q^2} \left\| \frac{a_{t|s} \cancel{\sigma_s^2}}{\cancel{\sigma_t^2}} z_t + \frac{a_s \sigma_{t|s}^2}{\sigma_t^2} x - \frac{a_{t|s} \cancel{\sigma_s^2}}{\cancel{\sigma_t^2}} z_t - \frac{a_s \sigma_{t|s}^2}{\sigma_t^2} \hat{n}_\theta \right\|_2^2 \\
 &= \frac{a_s^2 \sigma_{t|s}^4}{2 \sigma_Q^2 \sigma_t^4} \left\| x - \hat{n}_\theta \right\|_2^2 = \frac{a_s^2 (\sigma_t^2 - a_{t|s}^2 \sigma_s^2)}{2 \sigma_s^2 \sigma_t^2} \left\| x - \hat{n}_\theta \right\|_2^2 = \frac{1}{2} \left( \underbrace{\frac{a_s^2}{\sigma_s^2}}_{\text{SNR}(s)} - \underbrace{\frac{a_t^2}{\sigma_t^2}}_{\text{SNR}(t)} \right) \left\| x - \hat{n}_\theta \right\|_2^2
 \end{aligned}$$

$$\boxed{2.g} \quad \begin{aligned} s(i) &= \frac{i-1}{T} \\ t(i) &= \frac{i}{T} \end{aligned} \longrightarrow s = t - \frac{1}{T}$$

$$\begin{aligned}
 L_T(x) &= \frac{T}{2} \mathbb{E} \left[ \left( \text{SNR}\left(t - \frac{1}{T}\right) - \text{SNR}(t) \right) \left\| x - \hat{n}_\theta \right\|_2^2 \right] \\
 &= \frac{-1}{2} \mathbb{E} \left[ \frac{\text{SNR}\left(t - \frac{1}{T}\right) - \text{SNR}(t)}{-\frac{1}{T}} \left\| x - \hat{n}_\theta \right\|_2^2 \right]
 \end{aligned}$$

$$\lim_{T \rightarrow \infty} L_T(x) = -\frac{1}{2} \mathbb{E} \left[ \text{SNR}'(t) \left\| x - \hat{n}_\theta \right\|_2^2 \right]$$



3.1.a

اگر عبارت فوق را بنویسیم داریم:  $x_{t+1} = x_t + \delta \nabla_x \log p(x_t)$

همان طور که مشاهده می شود داریم بر روی عملیات  $p(x)$  gradient ascent را انجام می دهیم (در فضای  $x \in \mathbb{R}^d$  توزیع را بهینه می کنیم) بنابراین به یکی از قله های  $p(x)$  خواهیم رسید.

در MLE مشتق توزیع را باید حساب کنیم (نسبت به پارامترها) تا بتوانیم  $\hat{\theta}_{MLE}$  را به دست آوریم و به توزیع مد نظر برسیم. همچنین فرض MLE آن است که خانواده توزیع را می دانیم. به طور کلی ممکن است MLE جواب نداشته باشد زیرا ممکن است نتوانیم مشتق را به صورت نرم بسته به دست آوریم. در عبارت Langevin Dynamics مشتق را نسبت به  $x$  می گیریم نه  $\theta$ . و ضمناً یک فرایند iterative را حساب می کنیم و gradient ascent می زنیم تا به بهینه محلی برسیم.

3.1.b

دلیل استفاده از نویز ع. ایجاد randomness در فرایند generation و بنابراین ایجاد تنوع در سیمپل گیری است. به این صورت حتی اگر از نقطه شروع یکسانی نیز شروع به حرکت کنیم، هر دفعه به سیمپل های متفاوتی می برسیم.

اگر قله یافت شده نیز باشد باید به سببی که در راسهای مختلف آن نقطه وجود دارد ممکن است حرکتی انجام شود اما اگر flat باشد در آن نقطه گیر می کند.



$$\boxed{3.2.a} \quad \nabla_n \log \underbrace{\frac{1}{M} \sum_{i=1}^M \left( \frac{1}{\sqrt{2\pi} \sigma} \right)^d \exp\left( \frac{-\|n - x^{(i)}\|^2}{2\sigma^2} \right)}_{q(n)} = \frac{q'(n)}{q(n)}$$

$$= \frac{\frac{1}{M} \sum_{i=1}^M \left( \frac{1}{\sqrt{2\pi} \sigma} \right)^d \left[ -2 \times \frac{1}{2\sigma^2} (n - x^{(i)}) \exp\left( \frac{-\|n - x^{(i)}\|^2}{2\sigma^2} \right) \right]}{\frac{1}{M} \sum_{i=1}^M k(n | x^{(i)})}$$

$$= \frac{\sum_{i=1}^M \frac{1}{\sigma^2} (x^{(i)} - n) k(n | x^{(i)})}{\sum_{i=1}^M k(n | x^{(i)})}$$

$\boxed{3.2.b}$  در بخش قبل  $q(n)$  را با کرنل گاوسی تخمین زدیم (GMM) که لزوماً تخمین خوبی از توزیع نیست. به طور کلی در این روش به صورت explicit به دنبال یافتن توزیع هستیم اما روش implicit که در denoising score matching استفاده می‌شود بهتر است.

$\boxed{3.2.d}$

$$\nabla_n \log q(n | x_0) = \nabla_n \left[ -d \log(\sqrt{2\pi} \sigma) - \frac{\|n - x_0\|^2}{2\sigma^2} \right] = -\frac{1}{2\sigma^2} 2(n - x_0) = \frac{x_0 - n}{\sigma^2}$$

$$\rightarrow J_2(\theta) = \mathbb{E}_{q(n, x_0)} \left[ \frac{1}{2} \left\| s_\theta(n) + \frac{n - x_0}{\sigma^2} \right\|^2 \right]$$

می‌دانیم که در هر مرحله  $x$  نسبت به مراحل قبل دارای پیوستری دارد. یعنی نویز پیوستری به دقتور اضافه شده یا به اصطلاح manifold پیوسته یافت کرده. فرض می‌کنیم که داریم:  $q(x | x_0) = \mathcal{N}(x_0, \sigma^2 I)$

$$\rightarrow J_2(\theta) = \mathbb{E}_{\substack{q(n, x_0) \\ \varepsilon}} \left[ \frac{1}{2} \left\| \underbrace{s_\theta(n)}_{-\varepsilon_\theta(n)} + \frac{(x_0 + \sigma \varepsilon) - x_0}{\sigma^2} \right\|^2 \right] = \frac{1}{2\sigma^2} \mathbb{E}_{\substack{q(n_0) \\ \varepsilon}} \left[ \left\| \varepsilon - \varepsilon_\theta(x_0 + \sigma \varepsilon) \right\|^2 \right]$$

بنابراین در فرایند آموزش می‌فهمیم فاصله بین دینوایزر و نویز اصلی را کم کنیم. ضمناً در فرایند آموزش  $\mathbb{E}$  حذف شده و به ازای یک batch میانگین تجربی می‌گیریم.

برای sampling هم کافی است یک  $x_T \sim \mathcal{N}(0, I)$  را توسط فرمول  $x_{t-1} = x_t + \frac{\sigma}{2} s_\theta(x_t) + \sqrt{\sigma} z_t$  در طی  $T$  مرحله denoise کنیم  
 $z_t \sim \mathcal{N}(0, I)$

score matching  $\xrightarrow{(\mathcal{J}_1)}$   $\mathbb{E} \left[ \frac{1}{2} \| s_\theta(u) - \nabla_u \log q(u) \|^2 \right]$

3.2.c

$$= \underbrace{\mathbb{E} \left[ \frac{1}{2} \| s_\theta(u) \|^2 \right]}_{c_1} + \underbrace{\mathbb{E} \left[ \frac{1}{2} \| \nabla_u \log q(u) \|^2 \right]}_{c_2} - \underbrace{\mathbb{E} \left[ \langle s_\theta(u), \nabla_u \log q(u) \rangle \right]}_{g(\theta)}$$

$$g(\theta) = \int_u \cancel{q(u)} \langle s_\theta(u), \frac{\nabla_u q(u)}{\cancel{q(u)}} \rangle du = \int_u \langle s_\theta(u), \underbrace{\nabla_u \int_{n_0} q(u|n_0) q(n_0) dn_0}_{\int_{n_0} q(n_0) \underbrace{\nabla_u q(u|n_0)}_{q(u|n_0) \nabla_u \log q(u|n_0)} dn_0} \rangle du$$

$$= \int_u \int_{n_0} \overbrace{q(n_0) q(u|n_0)}^{q(u, n_0)} \langle s_\theta(u), \nabla_u \log q(u|n_0) \rangle dn_0 du$$

$$= \mathbb{E}_{q(u, n_0)} \left[ \langle s_\theta(u), \nabla_u \log q(u|n_0) \rangle \right]$$

$$\xrightarrow{(\mathcal{J}_1)} \text{score matching} = c_1 + c_2 + \overbrace{g(\theta)}^{\mathcal{J}_2} = \mathbb{E} \left[ \frac{1}{2} \| s_\theta(u) - \nabla_u \log q(u|n_0) \|^2 \right] - \mathbb{E} \left[ \frac{1}{2} \| \nabla_u \log q(u|n_0) \|^2 \right] + \mathbb{E} \left[ \frac{1}{2} \| \nabla_u \log q(u) \|^2 \right]$$

$$\xrightarrow{} \mathcal{J}_2(\theta) = \mathcal{J}_1(\theta) + \underbrace{\mathbb{E} \left[ \frac{1}{2} \| \nabla_u \log q(u|n_0) \|^2 \right] - \mathbb{E} \left[ \frac{1}{2} \| \nabla_u \log q(u) \|^2 \right]}_c$$

# 4.1.a

می دانیم که لاس برابر است با  $L(\theta, \psi) = f(\theta\psi) + c$  و بنابراین  $V(\theta, \psi) = \begin{bmatrix} \dot{\theta}(t) \\ \dot{\psi}(t) \end{bmatrix} = \begin{bmatrix} -f'(\theta\psi)\psi \\ f'(\theta\psi)\theta \end{bmatrix}$  دقت شود که مشتق نسبت به  $\theta$  علامت منفی دارد زیرا لاس generator منفی لاس discriminator است.

حال چون  $L(\theta, 0) = L(0, \psi)$  و برابر  $c$  است پس  $(\theta, \psi) = (0, 0)$  نقطه ای از نقاط تعادل است. حال می خواهیم نشان دهیم که این نقطه تنها نقطه ای تعادل است. فرض مسئله درباره  $f$  آن است که مشتق آن همواره مثبت است پس تنها حالتی که  $V(\theta, \psi)$  بتواند 0 شود آن است که جفت  $\theta, \psi$  صفر شوند و تنها نقطه تعادل  $(\theta, \psi) = (0, 0)$  است.

$$V'(\theta, \psi) = \begin{bmatrix} -f''(\theta\psi)\psi^2 & -f'(\theta\psi) - f''(\theta\psi)\theta\psi \\ f'(\theta\psi) + f''(\theta\psi)\theta\psi & f''(\theta\psi)\theta^2 \end{bmatrix}$$

$$\rightarrow V'(0, 0) = \begin{bmatrix} 0 & -f'(0) \\ f'(0) & 0 \end{bmatrix} \rightarrow \text{eig vals} = \pm f'(0)i$$

مشاهده می شود که هر دو مقدار ویژه مایوس را دارند  
معدومی است

# 4.1.b

اگر تعریف کنیم  $V_1(\theta, \psi) = \begin{bmatrix} -\nabla_{\theta} L(\theta, \psi) \\ 0 \end{bmatrix}$  و  $V_2(\theta, \psi) = \begin{bmatrix} 0 \\ \nabla_{\psi} L(\theta, \psi) \end{bmatrix}$  یعنی در صفت  $V$  را به  $V_1$  و  $V_2$  بسطیم،

$$\frac{d}{dt} \theta(t)^2 + \psi(t)^2 = 2\theta(t) V_1(\theta, \psi) + 2\psi(t) V_2(\theta, \psi) = 0 \quad \forall t \in [0, \infty)$$

بنابراین چون مشتق  $\theta^2 + \psi^2$  صفر شود پس خودی const است (به ازای هر  $t \in [0, \infty)$ )

# 4.1.c

طبق بحثی 4.1.a هر دو مقدار ویژه معدومی هستند پس به طور کلی درباره همگرایی نمی توان نظر داد ولی اگر همگرا باشد با نرخ زیر خطی همگرا است

# 4.2.a

$F_h(\theta, \psi) = (\theta, \psi) + hV(\theta, \psi)$  برای simultaneous gradient descent داریم

$$\rightarrow F'_h(\theta^*, \psi^*) = I + hV'(\theta^*, \psi^*)$$

بنابراین اگر eig val های  $V'(\theta^*, \psi^*)$  را  $\mu$  بنامیم  
آنگاه  $\lambda = 1 + h\mu$

$$\xrightarrow{\mu = -a + ib, a > 0} |\lambda|^2 = |1 + h(-a + ib)|^2 = (1 - ha)^2 + h^2 b^2 = (a^2 + b^2)h^2 - 2ah + 1$$

$$\xrightarrow{\text{if } |\lambda| > 1} (a^2 + b^2)h^2 - 2ah + 1 > 1 \rightarrow (a^2 + b^2)h - 2a > 0 \rightarrow h > \frac{2a}{a^2 + b^2}$$

بنابراین به ازای نرخ یادگیری ذکر شده، مقادیر ویژه بیرون دایره واحد تدرار می گیرند

$$(\theta_{k+1}, \psi_{k+1}) = (\theta_k, \psi_k) + h V(\theta_k, \psi_k) \quad \text{نقطة التالية}$$

$$\begin{aligned} \rightarrow \theta_{k+1}^2 + \psi_{k+1}^2 &= \left[ \theta_k + h V_1(\theta_k, \psi_k) \right]^2 + \left[ \psi_k + h V_2(\theta_k, \psi_k) \right]^2 \\ &= \theta_k^2 + h^2 V_1(\theta_k, \psi_k)^2 + \cancel{2 \theta_k h V_1(\theta_k, \psi_k)} + \psi_k^2 + h^2 V_2(\theta_k, \psi_k)^2 + \cancel{2 \psi_k h V_2(\theta_k, \psi_k)} \\ &\quad \quad \quad \underbrace{\hspace{10em}}_{-f'(\theta_k, \psi_k) \psi_k} \hspace{10em} \underbrace{\hspace{10em}}_{f'(\theta_k, \psi_k) \theta_k} \end{aligned}$$

بنابر این مشاهده می شود که تمامی عبارات توان 2 دارند و بنابر این تفاضل دو کلام بزرگتر مساوی صفر است. این کار مشابه مثبت شدن مشتق در حالت پیوسته است و نشان می دهد  $\theta^2 + \psi^2$  به صورت یک کلا افزایش می یابد.

:  $\rho_{\mu}$  alternating gradient descent  $\sigma_{\mu}$

update operators  $\rightarrow \begin{cases} F_1(\theta, \psi) = \begin{bmatrix} \theta - h f'(\theta\psi) \psi \\ \psi \end{bmatrix} \\ F_2(\theta, \psi) = \begin{bmatrix} \theta \\ \psi + h f'(\theta\psi) \theta \end{bmatrix} \end{cases}$

$$\rightarrow \begin{cases} F_1'(0,0) = \begin{bmatrix} 1 & -h f'(0) \\ 0 & 1 \end{bmatrix} \\ F_2'(0,0) = \begin{bmatrix} 1 & 0 \\ h f'(0) & 1 \end{bmatrix} \end{cases} \xrightarrow[\text{update operator}]{\text{Jacobian of the combined}} F_1'(0,0) \circ F_2'(0,0) = \begin{bmatrix} 1 & -h f'(0) \\ h f'(0) & 1 - h^2 f'(0)^2 \end{bmatrix}$$

$$\rightarrow \text{eig vals} = \lambda_{1/2} = 1 - \frac{(\hbar F(\omega))^2}{2} \pm \sqrt{\left(1 - \frac{(\hbar F(\omega))^2}{2}\right)^2 - 1}$$

برای قسمت 4.2.b مقادیر ویژه را به دست آوریم و مشاهده می‌کنیم که به ازای  $2 \leq h f(\phi)$  روی دایره واحد می‌ایستیم. در این صورت در مورد همگرایی نمی‌توان نظر داد اما اگر  $h f(\phi) < 2$  باشد، با نرخ زیر خطی همگراست.

برای قسمت 4.2.a نیز در حالتی که مد نظر سوال است هر 2 مقدار دیرینه خارج دایره داده قرار دارند و بنابراین الگوریتم همگرا نیست

در صورت استفاده از نویز تابع هدف GAN به صورت زیر می شود:

4.3.a

$$\mathbb{E}_{\tilde{\theta} \sim \mathcal{N}(\theta, \epsilon^2)} [f(\tilde{\theta}\psi)] + \mathbb{E}_{\kappa \sim \mathcal{N}(0, \epsilon^2)} [f(-\kappa\psi)]$$

$$\rightarrow \tilde{V}(\theta, \psi) = \mathbb{E}_{\tilde{\theta}, \kappa} \begin{bmatrix} -\psi f'(\tilde{\theta}\psi) \\ \tilde{\theta} f'(\tilde{\theta}\psi) - \kappa f'(-\kappa\psi) \end{bmatrix}$$

$$\rightarrow \tilde{V}'(\theta, \psi) = \mathbb{E}_{\tilde{\theta}, \kappa} \begin{bmatrix} -f''(\tilde{\theta}\psi)\psi^2 & -f'(\tilde{\theta}\psi) - f''(\tilde{\theta}\psi)\tilde{\theta}\psi \\ f'(\tilde{\theta}\psi) + f''(\tilde{\theta}\psi)\tilde{\theta}\psi & f''(\tilde{\theta}\psi)\tilde{\theta}^2 + \kappa^2 f''(-\kappa\psi) \end{bmatrix}$$

$$\rightarrow \tilde{V}'(0, 0) = \begin{bmatrix} 0 & -f'(0) \\ f'(0) & 2f''(0)\epsilon^2 \end{bmatrix} \rightarrow \text{eig vals} = \lambda_{\frac{1}{2}} = f''(0)\epsilon^2 \pm \sqrt{f''(0)^2\epsilon^4 - f'(0)^2}$$

4.3.b

$$f'(t) = -\frac{-e^{-t}}{1+e^{-t}} \rightarrow f'(0) = \frac{1}{2}$$

برای حالت پیوسته در بخش 4.1.a نشان دادیم که مقادیر ویژه  $\pm f'(0)i$  هستند  
بنابراین به ازای  $f(t)$  داده شده مقادیر ویژه  $\pm \frac{1}{2}i$  می شوند که با فرض زیر قطبی هستند

در برخی سوال ها با علی بابایی همکاری داشتیم