

CS 236 Autumn 2019/2020 Homework 3

SUNet ID: 06009508

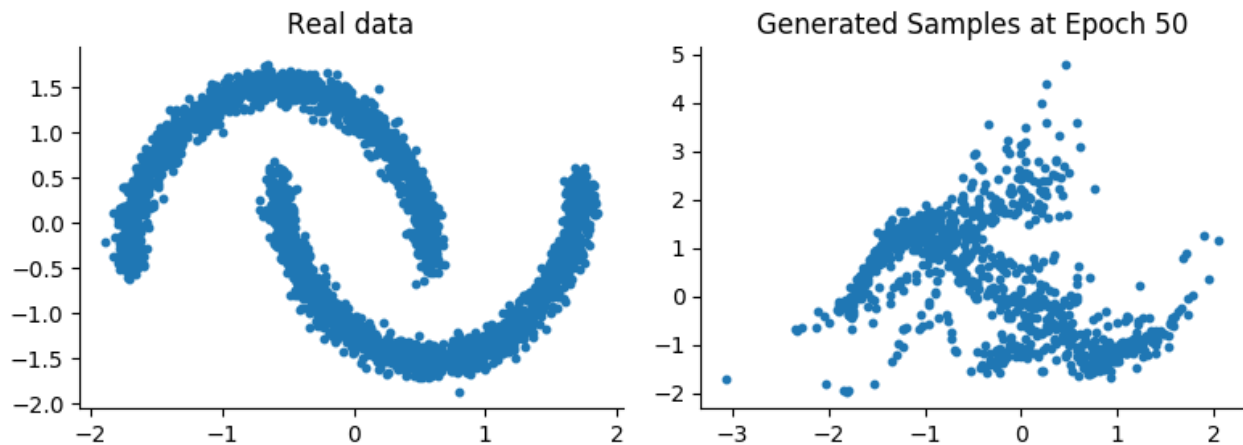
Name: Brandon McKinzie

Collaborators: N/A

By turning in this assignment, I agree by the Stanford honor code and declare that all of this is my own work.

Problem 2: Generative adversarial networks (15 points)

Part 4.



Problem 2: Generative adversarial networks (15 points)

Unfortunately, this form of loss for L_G suffers from a vanishing gradient problem. In terms of the discriminator's logits, the minimax loss is

$$L_G^{\text{minimax}}(\theta; \phi) = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)} [\log(1 - \sigma(h_\phi(G_\theta(\mathbf{z})))]$$

Show that the derivative of L_G^{minimax} with respect to θ is approximately 0 if $D(G_\theta(\mathbf{z})) \approx 0$, or equivalently, if $h_\phi(G_\theta(\mathbf{z})) \ll 0$. You may use the fact that $\sigma'(x) = \sigma(x)(1 - \sigma(x))$. Why is this problematic for the training of the generator when the discriminator successfully identifies a fake sample $G_\theta(\mathbf{z})$?

Below, I take the derivative with respect to θ and evaluate it for the case of $\sigma(h_\phi(G_\theta(\mathbf{z}))) \approx 0$.

$$\frac{\partial L_G^{\text{minimax}}(\theta; \phi)}{\partial \theta} = \frac{\partial}{\partial \theta} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)} [\log(1 - \sigma(h_\phi(G_\theta(\mathbf{z})))]) \quad (1)$$

$$= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)} \left[\frac{\partial}{\partial \theta} \log(1 - \sigma(h_\phi(G_\theta(\mathbf{z}))) \right] \quad (2)$$

$$\frac{\partial}{\partial \theta} \log(1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) = \frac{1}{1 - \sigma(h_\phi(G_\theta(\mathbf{z})))} \frac{\partial}{\partial \theta} (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) \quad (3)$$

$$= -\frac{1}{1 - \sigma(h_\phi(G_\theta(\mathbf{z})))} \frac{\partial}{\partial \theta} \sigma(h_\phi(G_\theta(\mathbf{z}))) \quad (4)$$

$$\frac{\partial}{\partial \theta} \sigma(h_\phi(G_\theta(\mathbf{z}))) = \sigma(h_\phi(G_\theta(\mathbf{z}))) (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) \frac{\partial}{\partial \theta} h_\phi(G_\theta(\mathbf{z})) \quad (5)$$

$$\therefore \left. \frac{\partial L_G^{\text{minimax}}(\theta; \phi)}{\partial \theta} \right|_{h_\phi < 0} = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)} \left[-\frac{\sigma(h_\phi(G_\theta(\mathbf{z}))) (1 - \sigma(h_\phi(G_\theta(\mathbf{z})))) \frac{\partial}{\partial \theta} h_\phi(G_\theta(\mathbf{z}))}{1 - \sigma(h_\phi(G_\theta(\mathbf{z})))} \right] \quad (6)$$

$$= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)} \left[-\sigma(h_\phi(G_\theta(\mathbf{z}))) \frac{\partial}{\partial \theta} h_\phi(G_\theta(\mathbf{z})) \right] \quad (7)$$

$$\approx 0 \quad (8)$$

This is problematic since it suggests that if our generator G_θ is horrible – that is, it produces obviously fake samples from the perspective of our discriminator (or, equivalently, if our discriminator is perfect) – then its gradients will be (nearly) zero and thus its parameters won't get updated.



Problem 3: Divergence minimization (25 points)

Part 1. Show that L_D is minimized when $D_\phi = D^*$, where

$$D^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})}.$$

(Hint: for a fixed \mathbf{x} , what t minimizes $f(t) = -p_{data}(\mathbf{x}) \log t - p_\theta(\mathbf{x}) \log(1 - t)$?)

First we expand out the equation for L_D as an integral:

$$L_D(\phi; \theta) = - \int [p_{data}(\mathbf{x}) \log D_\phi(\mathbf{x}) + p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))] d\mathbf{x} \quad (9)$$

The derivative of the integrand with respect to D is

$$p_{data}(\mathbf{x}) \frac{1}{D_\phi(\mathbf{x})} - p_\theta(\mathbf{x}) \frac{1}{1 - D_\phi(\mathbf{x})} \quad (10)$$

We can set this to zero and solve for D to obtain D^* .

$$p_{data}(\mathbf{x}) \frac{1}{D_\phi(\mathbf{x})} = p_\theta(\mathbf{x}) \frac{1}{1 - D_\phi(\mathbf{x})} \quad (11)$$

$$p_{data}(\mathbf{x})(1 - D_\phi(\mathbf{x})) = p_\theta(\mathbf{x}) D_\phi(\mathbf{x}) \quad (12)$$

$$p_{data}(\mathbf{x}) = D_\phi(\mathbf{x})(p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})) \quad (13)$$

$$D_\phi(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})} \quad (14)$$

Part 2. Recall that $D_\phi(\mathbf{x}) = \sigma(h_\phi(\mathbf{x}))$. Show that the logits $h_\phi(\mathbf{x})$ of the discriminator estimate the log of the likelihood ratio of \mathbf{x} under the true distribution compared to the model's distribution; that is, show that if $D_\phi = D^*$, then

$$h_\phi(\mathbf{x}) = \log \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x})}$$

We just plug in $D_\phi(\mathbf{x}) = \sigma(h_\phi(\mathbf{x}))$ and solve for $h_\phi(\mathbf{x})$ using basic arithmetic:

$$\sigma(h_\phi(\mathbf{x})) = \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \quad (15)$$

$$\text{Let } F := \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \quad (16)$$

$$\frac{1}{1 + e^{-h_\phi(\mathbf{x})}} = F \quad (17)$$

$$1 = (1 + e^{-h_\phi(\mathbf{x})})F \quad (18)$$

$$= F + Fe^{-h_\phi(\mathbf{x})} \quad (19)$$

$$\frac{1 - F}{F} = e^{-h_\phi(\mathbf{x})} \quad (20)$$

$$\log \frac{1 - F}{F} = -h_\phi(\mathbf{x}) \quad (21)$$

$$\therefore h_\phi(\mathbf{x}) = \log \frac{F}{1 - F} \quad (22)$$

$$= \log \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x})} \quad (23)$$

Part 3. Consider a generator loss defined by the sum of the minimax loss and the non-saturating loss,

$$L_G(\theta; \phi) = \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[\log(1 - D_\phi(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[\log D_\phi(\mathbf{x})].$$

Show that if $D_\phi = D^*$, then

$$L_G(\theta; \phi) = KL(p_\theta(\mathbf{x}) || p_{data}(\mathbf{x})).$$

Again, we just plug in for D^* and shuffle some terms around to arrive at the result.

$$L_G(\theta; \phi) = \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[\log(\frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})})] - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[\log \frac{p_{data}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})}] \quad (24)$$

$$= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[\log \frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})} / \frac{p_{data}(\mathbf{x})}{p_\theta(\mathbf{x}) + p_{data}(\mathbf{x})}] \quad (25)$$

$$= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[\log \frac{p_\theta(\mathbf{x})}{p_{data}(\mathbf{x})}] \quad (26)$$

$$= KL(p_\theta(\mathbf{x}) || p_{data}(\mathbf{x})) \quad (27)$$

Part 4. Recall that when training VAEs, we minimize the negative ELBO, an upper bound to the negative log likelihood. Show that the negative log likelihood, $-\mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})]$, can be written as a KL divergence plus an additional term that is constant with respect to θ .

We can use the definition of $KL(p_{data}(x)||p_{\theta}(x))$ and rearrange some terms:

$$KL(p_{data}(x)||p_{\theta}(x)) \triangleq \mathbb{E}_{x \sim p_{data}(x)} \left[\log \frac{p_{data}(x)}{p_{\theta}(x)} \right] \quad (28)$$

$$= \mathbb{E}_{x \sim p_{data}(x)} [\log p_{data}(x)] - \mathbb{E}_{x \sim p_{data}(x)} [\log p_{\theta}(x)] \quad (29)$$

$$\therefore -\mathbb{E}_{x \sim p_{data}(x)} [\log p_{\theta}(x)] = KL(p_{data}(x)||p_{\theta}(x)) - \mathbb{E}_{x \sim p_{data}(x)} [\log p_{data}(x)] \quad (30)$$

and we see that $-\mathbb{E}_{x \sim p_{data}(x)} [\log p_{\theta}(x)]$ can be written as $KL(p_{data}(x)||p_{\theta}(x))$ plus a θ -independent term, $-\mathbb{E}_{x \sim p_{data}(x)} [\log p_{data}(x)]$. This of course implies that

$$\min_{\theta} -\mathbb{E}_{x \sim p_{data}(x)} [\log p_{\theta}(x)] = \min_{\theta} \left[KL(p_{data}(x)||p_{\theta}(x)) - \mathbb{E}_{x \sim p_{data}(x)} [\log p_{data}(x)] \right] \quad (31)$$

$$= \min_{\theta} KL(p_{data}(x)||p_{\theta}(x)) \quad (32)$$

Does this mean that a VAE decoder trained with ELBO and a GAN generator trained with the L_G defined in the previous part are implicitly learning the same objective? Explain.

No, they are not the same, since the GAN generator from the previous part minimizes

$$L_G(\theta; \phi) = KL(p_{\theta}(x)||p_{data}(x)) \quad (33)$$

and $KL(p_{\theta}(x)||p_{data}(x)) \neq KL(p_{data}(x)||p_{\theta}(x))$ since the KL divergence is not symmetric.

Problem 4: Conditional GAN with projection discriminator (20 points)

Part 1. Suppose that when $(\mathbf{x}, y) \sim p_{data}(\mathbf{x}, y)$, there exists a feature mapping φ under which $\varphi(\mathbf{x})$ becomes a mixture of m unit Gaussians, with one Gaussian per class label y . Assume that when $(\mathbf{x}, y) \sim p_{\theta}(\mathbf{x}, y)$, $\varphi(\mathbf{x})$ also becomes a mixture of m unit Gaussians, again with one Gaussian per class label y . Concretely, we assume that the ratio of the conditional probabilities can be written as

$$\frac{p_{data}(\mathbf{x}|y)}{p_{\theta}(\mathbf{x}|y)} = \frac{\mathcal{N}(\varphi(\mathbf{x})|\boldsymbol{\mu}_y, I)}{\mathcal{N}(\varphi(\mathbf{x})|\hat{\boldsymbol{\mu}}_y, I)},$$

where $\boldsymbol{\mu}_y$ and $\hat{\boldsymbol{\mu}}_y$ are the means of the Gaussians for p_{data} and p_{θ} respectively.

Show that under this simplifying assumption, the optimal discriminator's logits $h^*(\mathbf{x}, y)$ can be written in the form

$$h^*(\mathbf{x}, y) = \mathbf{y}^T (A\varphi(\mathbf{x}) + \mathbf{b})$$

for some matrix A and vector \mathbf{b} , where \mathbf{y} is a one-hot vector denoting the class y . In this problem, the discriminator's output and logits are related by $D_{\phi}(\mathbf{x}, y) = \sigma(h_{\phi}(\mathbf{x}, y))$. (Hint: use the result from problem 2.2.)

To rephrase the question a bit for clarity: there exists some mapping $\varphi : (\mathbf{x} \in \mathbb{R}^d, y) \mapsto \mathbf{r} \in \mathbb{R}^r$ where \mathbf{r} is modeled by the distribution $\mathcal{N}(\mathbf{r}; \boldsymbol{\mu}_y, I)$. The result from problem 3.2 states that the optimal unconditional discriminator's logits $h^*(\mathbf{x})$ satisfy

$$h^*(\mathbf{x}) = \log \frac{p_{data}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \quad (34)$$

For the class-conditional setting, our optimal discriminator logits $h^*(\mathbf{x}, y)$ will instead follow

$$h^*(\mathbf{x}, y) = \log \frac{p_{data}(\mathbf{x} | y)}{p_{\theta}(\mathbf{x} | y)} \quad (35)$$

$$= \log \frac{\mathcal{N}(\varphi(\mathbf{x})|\boldsymbol{\mu}_y, I)}{\mathcal{N}(\varphi(\mathbf{x})|\hat{\boldsymbol{\mu}}_y, I)} \quad (36)$$

$$= \log \frac{\exp\left(-\frac{1}{2}(\varphi(\mathbf{x}) - \boldsymbol{\mu}_y)^T(\varphi(\mathbf{x}) - \boldsymbol{\mu}_y)\right)}{\exp\left(-\frac{1}{2}(\varphi(\mathbf{x}) - \hat{\boldsymbol{\mu}}_y)^T(\varphi(\mathbf{x}) - \hat{\boldsymbol{\mu}}_y)\right)} \quad (37)$$

$$= -\frac{1}{2}(\varphi(\mathbf{x}) - \boldsymbol{\mu}_y)^T(\varphi(\mathbf{x}) - \boldsymbol{\mu}_y) + \frac{1}{2}(\varphi(\mathbf{x}) - \hat{\boldsymbol{\mu}}_y)^T(\varphi(\mathbf{x}) - \hat{\boldsymbol{\mu}}_y) \quad (38)$$

$$= \varphi(\mathbf{x})^T [\boldsymbol{\mu}_y - \hat{\boldsymbol{\mu}}_y] + \frac{1}{2}(\hat{\boldsymbol{\mu}}_y^2 - \boldsymbol{\mu}_y^2) \quad (39)$$

We can equivalently rewrite this to use the one-hot vector \mathbf{y} by defining an m -element vector \mathbf{b} whose i th entry equals $\frac{1}{2}(\hat{\boldsymbol{\mu}}_i^2 - \boldsymbol{\mu}_i^2)$ (corresponding to class i). Similarly, we define the matrix \mathbf{A} with m rows and r columns (dimensionality of $\varphi(\mathbf{x})$), with $A_{i,:} := \boldsymbol{\mu}_i - \hat{\boldsymbol{\mu}}_i$. This results in

$$h^*(\mathbf{x}, y) = (\mathbf{A}\varphi(\mathbf{x}) + \mathbf{b})_y \quad (40)$$

$$= \mathbf{y}^T (\mathbf{A}\varphi(\mathbf{x}) + \mathbf{b}) \quad (41)$$



Problem 5: Wasserstein GAN (35 points)

Part 1. Let $p_\theta(x) = \mathcal{N}(x|\theta, \epsilon^2)$ and $p_{data}(x) = \mathcal{N}(x|\theta_0, \epsilon^2)$ be normal distributions with standard deviation ϵ centered at $\theta \in \mathbb{R}$ and $\theta_0 \in \mathbb{R}$ respectively. Show that

$$\text{KL}(p_\theta(x)||p_{data}(x)) = \frac{(\theta - \theta_0)^2}{2\epsilon^2}.$$

$$\text{KL}(p_\theta(x)||p_{data}(x)) = \mathbb{E}_{x \sim p_\theta} \left[\log \frac{p_\theta(x)}{p_{data}(x)} \right] \quad (42)$$

$$= \int_{\mathbb{R}} p_\theta(x) \log \frac{p_\theta(x)}{p_{data}(x)} dx \quad (43)$$

$$= \int_{\mathbb{R}} \mathcal{N}(x | \theta, \epsilon^2) \log \frac{\exp\left(-\frac{(x-\theta)^2}{2\epsilon^2}\right)}{\exp\left(-\frac{(x-\theta_0)^2}{2\epsilon^2}\right)} dx \quad (44)$$

$$= \int_{\mathbb{R}} \mathcal{N}(x | \theta, \epsilon^2) \left(-\frac{1}{2\epsilon^2} \left((x-\theta)^2 - (x-\theta_0)^2 \right) \right) dx \quad (45)$$

$$= \int_{\mathbb{R}} \mathcal{N}(x | \theta, \epsilon^2) \left(-\frac{1}{2\epsilon^2} \left(-2x\theta + \theta^2 + 2x\theta_0 - \theta_0^2 \right) \right) dx \quad (46)$$

$$= \frac{1}{2\epsilon^2} \int_{\mathbb{R}} \mathcal{N}(x | \theta, \epsilon^2) \left(2x(\theta - \theta_0) + \theta_0^2 - \theta^2 \right) dx \quad (47)$$

$$= \frac{1}{2\epsilon^2} \mathbb{E}_{x \sim p_\theta} \left[2x(\theta - \theta_0) + \theta_0^2 - \theta^2 \right] \quad (48)$$

$$= \frac{1}{2\epsilon^2} \left(2\theta(\theta - \theta_0) + \theta_0^2 - \theta^2 \right) \quad (49)$$

$$= \frac{1}{2\epsilon^2} \left(2\theta^2 - 2\theta\theta_0 + \theta_0^2 - \theta^2 \right) \quad (50)$$

$$= \frac{1}{2\epsilon^2} \left(\theta^2 - 2\theta\theta_0 + \theta_0^2 \right) \quad (51)$$

$$= \frac{(\theta - \theta_0)^2}{2\epsilon^2} \quad (52)$$

Part 2. Suppose $p_\theta(x)$ and $p_{data}(x)$ both place probability mass in only a very small part of the domain; that is, consider the limit $\epsilon \rightarrow 0$. What happens to $KL(p_\theta(x)||p_{data}(x))$ and its derivative with respect to θ , assuming that $\theta \neq \theta_0$? Why is this problematic for a GAN trained with the loss function L_G defined in problem 2.3?

As $\epsilon \rightarrow 0$, $KL \rightarrow \infty$. It's derivative,

$$\frac{\partial}{\partial \theta} \frac{(\theta - \theta_0)^2}{2\epsilon^2} = \frac{2(\theta - \theta_0)}{2\epsilon^2} = \frac{\theta - \theta_0}{\epsilon^2} \quad (53)$$

also goes to ∞ as $\epsilon \rightarrow 0$. We saw in problem 3.3 that if $D_\phi = D^*$, then $L_G(\theta, \phi) = KL(p_\theta(x)||p_{data}(x))$. Therefore, if we have a sufficiently good discriminator $\approx D^*$, small $\epsilon \rightarrow 0$, and $\theta \neq \theta_0$, our losses and gradients will explode and prevent our model from training successfully.

Part 3. To avoid this problem, we'll propose an alternative objective for the discriminator and generator. Consider the following alternative objectives:

$$\begin{aligned} L_D(\phi; \theta) &= \mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)] - \mathbb{E}_{x \sim p_{data}(x)}[D_\phi(x)] \\ L_G(\theta; \phi) &= -\mathbb{E}_{x \sim p_\theta(x)}[D_\phi(x)], \end{aligned}$$

where D_ϕ is no longer constrained to functions that output a probability; instead D_ϕ can be a function that outputs any real number. As defined, however, these losses are still problematic. Again consider the limit $\epsilon \rightarrow 0$; that is, let $p_\theta(x)$ be the distribution that outputs $\theta \in \mathbb{R}$ with probability 1, and let $p_{data}(x)$ be the distribution that outputs $\theta_0 \in \mathbb{R}$ with probability 1. Why is there no discriminator D_ϕ that minimizes this new objective L_D ?

In this case, L_D simply becomes

$$L_D(\phi, \theta) = D_\phi(\theta) - D_\phi(\theta_0) \quad (54)$$

Proceed via proof by contradiction. Assume that there exists some D^* that minimizes L_D . Since D can be any function that outputs a real number, let another function D^{**} be defined as

$$D^{**} = \begin{cases} D^*(\theta) - 1 & \text{if } x = \theta \\ D^*(\theta_0) & \text{if } x = \theta_0 \\ 0 & \text{otherwise} \end{cases} \quad (55)$$

Plugging this into L_D reveals

$$L_{D^{**}}(\phi, \theta) = D^{**}(\theta) - D^{**}(\theta_0) \quad (56)$$

$$= D^*(\theta) - 1 - D^*(\theta_0) \quad (57)$$

$$= L_{D^*}(\phi, \theta) - 1 \quad (58)$$

and we see that $L_{D^{**}} < L_{D^*}$, which contradicts our premise that $D^* = \arg \min_D L_D(\phi, \theta)$. Therefore, there does not exist any function D that minimizes this form of L_D (under the assumptions of p_θ , p_{data} , etc).

Part 4. *Let's tweak the alternate objective so that an optimal discriminator exists. Consider the same objective L_D and the same limit $\epsilon \rightarrow 0$. Now, suppose that D_ϕ is restricted to differentiable functions whose derivative is always between -1 and 1 . It can still output any real number. Is there now a discriminator D_ϕ out of this class of functions that minimizes L_D ? Briefly describe what the optimal D_ϕ looks like as a function of x .*

If we add the constraint that $-1 < \frac{\partial D_\phi(x)}{\partial x} < 1$, then

$$-|\theta - \theta_0| < L_D(\phi, \theta) = D_\phi(\theta) - D_\phi(\theta_0) < |\theta - \theta_0| \quad (59)$$

$$\inf L_D(\phi, \theta) = -|\theta - \theta_0| \quad (60)$$

If, for example, we had $\theta \in \mathbb{R}^+$ only, then $D_\phi(x) = -x$ would satisfy this.

