

## **Summary and Analysis of: ‘Fighting Gradients with Gradients: Dynamic Defenses against Adversarial Attacks’**

As we have learned throughout this project, deep networks can be extremely vulnerable to what are called adversarial attacks. This is where an adversary utilizes an imperceptible amount of noise. More technically a perturbation constrained by some norm threshold. This noise fools the classifier into incorrectly classify the given input. However, many of our current defenses are what is called “static”, meaning they do not change after training time. This leads to a lack of adaptivity for our deep network leads to trouble if the adversarial attacks evolve. This paper proposes a solution to this problem in a “dynamic” defense, meaning it updates during testing, called *Defensive entropy minimization*, or dent. What this means, is that any adversary never faces the same defense, it is constantly changing, making it more effective. The reason being, an adversary can find a hole in the static defense and continually exploit it and our defense can do nothing to prevent this. The paper discusses three main contributions, the weakness of static vs strengths of dynamic approaches, the specifics of dent – the dynamic model proposed, and the effectiveness of dent. The basic idea of the paper is most effectively put on page 3 at the start of section 3 “Adversarial attacks optimize against defenses at test time, so defenses should fight back, and counter-optimize against attacks.” Dent not only evolves during test-time, it only evolves during test-time, and does not change the training in any way. So while the adversary is optimizing the perturbation  $\delta$ , dent is optimizing both its model ( $\Delta$ ) and the input ( $\Sigma$ ), which both depend on the testing data. It updates on batches based on:

$$\text{argmin}_{\Sigma, \Delta} H(f(g(x + \delta; \Sigma); \theta + \Delta))$$

Where  $H$  is the Shannon Entropy:

$$H(y_i) = \sum_i -p(y_i) \log p(y_i)$$

It allows its defense parameters to reset between the batches. This allows the defense to move with the adversary as an attack of  $x + \delta^t$  can be met with  $\Delta^t$  which means the defense will always get the last move, which is an advantage. The model and input parameters are differentiable which is how dent ‘fights gradients with gradients.’ They tested a static method versus dent versus dent + (which is dent but adapting

sample-wise rather than batch-wise). We find that dent and dent + are way more effective in a basic example than the static method at improving even state-of-the-art defenses against adversaries or models with no defense at all, while still maintaining natural (no adversary) accuracy. Maintaining natural accuracy is extremely important, as often as a side effect, adversarial defenses can negatively impact natural accuracy. Dent is also robust to many different attack types, improving the best defenses by up to 20 points against certain attacks. Additionally, it is effective within many different architectures and different types of datasets. Ultimately although Dent has some limitations, such as depending on batches to evolve since it utilizes the entropy function as its update function, it seems to be way more effective and robust than even the best static defense in a big variety of situations.

4b.  $g(\vec{v}_2, \vec{v}_3) = \min_{\vec{z}} L(\vec{z}, \vec{v})$

$$= \min_{\vec{z}} \vec{c}^T \vec{z}_3 + \vec{v}_3^T \vec{z}_3 + \mathbf{1}_{S \in (\vec{x})}(\vec{z}_1) - \vec{v}_2^T W_1 \vec{z}_1$$

$$+ \left( \sum_{i=1}^n \mathbf{1}_{z_i}(\vec{z}_{2i}, \vec{z}_{2i}) - \vec{v}_3^T W_2 \vec{z}_2 + \vec{v}_2^T \vec{z}_2 \right) - \sum_{i=1}^n \vec{v}_{3i}^T b_i$$

since these terms are all only in terms of

at most one  $\vec{z}_i$  we can split our minimization

$$= \min_{\vec{z}_3} ((\vec{c} + \vec{v}_3)^T \vec{z}_3) + \min_{\vec{z}_1} (\mathbf{1}_{S \in (\vec{x})}(\vec{z}_1) - \vec{v}_2^T W_1 \vec{z}_1)$$

$$+ \left( \sum_{i=1}^n \left( \min_{\vec{z}_{2i}, \vec{z}_{2i}} (\mathbf{1}_{z_i}(\vec{z}_{2i}, \vec{z}_{2i}) - \vec{v}_3^T (W_2)_i \vec{z}_{2i} + \vec{v}_{2i}^T \vec{z}_{2i}) \right) - \underbrace{\sum_{i=1}^n \vec{v}_{3i}^T b_i}_{\text{not depend on } \vec{z}} \right)$$

As we wanted to show

4c.  $d^*(\vec{x}, \vec{c}) = \max_{\vec{v}} g(\vec{v}_2, \vec{v}_3)$

so we can first find the minima of each

$$\frac{\partial}{\partial z_3} ((\vec{c} + \vec{v}_3)^T \vec{z}_3) = 0$$

$$\frac{\partial}{\partial z_1} (\mathbf{1}_{S \in (\vec{x})}(\vec{z}_1) - \vec{v}_2^T W_1 \vec{z}_1) = 0$$

$$\mathbf{1}_{S \in (\vec{x})}(-\vec{v}_1^T \vec{v}_2) = 0$$

which will happen at the optimum of our  $\mathbf{1}_{S \in (\vec{x})}$

$$-\mathbf{1}_{S \in (\vec{x})}^*(W_1 \vec{v}_2)$$

$$\frac{\partial}{\partial z_1} \left( \mathbf{1}_{z_i}(\vec{z}_{2i}, \vec{z}_{2i}) - \vec{v}_3^T (W_2)_i \vec{z}_{2i} + \vec{v}_{2i}^T \vec{z}_{2i} \right) = 0$$

$$\sum_{i=1}^n \mathbf{1}_{z_i}(-\vec{v}_3^T (W_2)_i + \vec{v}_{2i}) = 0$$

similarly at optimum at character set etc  
function this is 0

$$\sum_{i=1}^n -\mathbf{1}_{z_i}^*(\vec{v}_3^T (W_2)_i - \vec{v}_{2i})$$

And everything else is fixed so we have

$$d^*(\vec{x}, \vec{c}) = \max_{\vec{v}} -\mathbf{1}_{S \in (\vec{x})}^*(W_1 \vec{v}_2) + \sum_{i=1}^n -\mathbf{1}_{z_i}^*(\vec{v}_3^T (W_2)_i - \vec{v}_{2i}) - \sum_{i=1}^n \vec{v}_{3i}^T b_i$$

s.t.  $\vec{v}_3 = -\vec{c}$  which is what we wanted to show

And using the definition as given we have:

$$d^*(\vec{x}, \vec{c}) = \max_{\vec{v}} -\mathbf{1}_{S \in (\vec{x})}^*(\vec{v}_1) + \sum_{i=1}^n -\mathbf{1}_{z_i}^*(\vec{v}_{2i} - \vec{v}_{3i}) - \sum_{i=1}^n \vec{v}_{3i}^T b_i$$

$$\text{s.t. } \vec{v}_3 = -\vec{c}$$

$$\vec{v}_2 = W_2^T \vec{v}_3$$

$$\vec{v}_1 = W_1^T \vec{v}_2$$

9.

$$\max_{\|x-x'\|_1 \leq \epsilon} L(F_\theta(x'), y_{\text{true}}) = \max_{z_n \in Z_\theta(x)} L(\hat{z}_n, y)$$

Since  $L(x, y) \leq L(x - a\mathbf{1}, y) + a$

$$\begin{aligned} \max_{\hat{z}_n \in Z_\theta(x)} L(\hat{z}_n, y) &\leq \max_{\hat{z}_n \in Z_\theta(x)} L(\hat{z}_n - (z_n)_* \mathbf{1}, y) \\ &= \max_{\hat{z}_n \in Z_\theta(x)} L((I - c_y \mathbf{1}^\top) z_n, y) \\ &= \max_{\hat{z}_n \in Z_\theta(x)} L(C \hat{z}_n, y) \end{aligned}$$

Letting  $C$  be the matrix  $(I - c_y \mathbf{1}^\top)$

by bounding the elementwise maximum over  $\hat{z}_n$  and the elementwise minimum for  $\bar{c} = y$

$$\max_{\hat{z}_n \in Z_\theta(x)} L(C \hat{z}_n, y) \leq L(h(\hat{z}_n))$$

Where  $\bar{c}$  being representative of the rows of  $C$  we define it as

$$h(z_k)_i = \max_{\hat{z}_n \in Z_\theta(x)} C_i \hat{z}_n$$

which gives us that

$$\begin{aligned} J_C(x, g_\theta(-c_i)) &\leq \min_{\hat{z}_n \in Z_\theta(x)} C_i^\top \hat{z}_n \\ -J_C(x, g_\theta(-c_i)) &\geq \max_{\hat{z}_n \in Z_\theta(x)} C_i^\top \hat{z}_n \\ h(z_k)_i &\leq -J_C(x, g_\theta(-c_i)) \end{aligned}$$

This all gives us that indeed

$$\max_{\|x-x'\|_1 \leq \epsilon} L(F_\theta(x'), y) \leq L(-J_C(x, g_\theta(c_y \mathbf{1}^\top - I)), y)$$

$$\begin{aligned}
 1. \quad & \underset{\vec{x}'}{\text{argmax}} \quad L(f_{\theta}(\vec{x}), \vec{y}_{\text{true}}) + \nabla_{\vec{x}} L(f_{\theta}(\vec{x}), y_{\text{true}})^T \vec{x}' \\
 & \text{s.t. } \|\vec{x} - \vec{x}'\|_{\infty} \leq \epsilon \\
 & \text{Let } \vec{x}' = \vec{x} + \epsilon \vec{u} \rightarrow \\
 & \underset{\substack{\vec{x} \\ \vec{u}}}{\text{argmax}} \quad \nabla_{\vec{x}} L(f_{\theta}(\vec{x}), \vec{y}_{\text{true}})^T (\vec{x} + \epsilon \vec{u}) \\
 & \text{s.t. } \|\epsilon \vec{u}\| \leq \epsilon \rightarrow \|\vec{u}\| \leq 1 \\
 & = \underset{\vec{u}}{\text{argmax}} \quad \epsilon \nabla_{\vec{x}} L(f_{\theta}(\vec{x}), \vec{y}_{\text{true}})^T \vec{u} \\
 & \text{s.t. } \|\vec{u}\|_{\infty} \leq 1
 \end{aligned}$$

The max is  $\epsilon \|\nabla_{\vec{x}} L(f_{\theta}(\vec{x}), \vec{y}_{\text{true}})\|_1$ ,  
thus  $\vec{u}$  must be

$$\begin{aligned}
 & \text{sgn}(\nabla_{\vec{x}} L(f_{\theta}(\vec{x}), y_{\text{true}})), \text{ thus} \\
 & \underset{\vec{x}'}{\text{argmax}} = \vec{x} + \epsilon \text{sgn}(\nabla_{\vec{x}} L(f_{\theta}(\vec{x}), y_{\text{true}}))
 \end{aligned}$$

$$2a) f^*(y) = \sup_s \{ y s - |s| \}$$

We can partition into casework.

Case 1:  $y < -1$ . Then  $g(s) = y s - |s|$ ,

$$\begin{cases} g(s) \rightarrow -\infty & s > 0 \\ g(s) \rightarrow \infty & s < 0 \end{cases} \quad \text{Thus } f^*(y) \text{ st. } y < -1 = \infty$$

Case 2:  $-1 < y < 1$ . Then  $g(s)$  follows the graph of  $-|s|$  closely, thus

$$\begin{cases} g(s) \rightarrow -\infty & s > 0 \\ g(s) \rightarrow -\infty & s < 0 \end{cases} \quad \begin{matrix} \text{with an apex} \\ \text{at } (0, 0). \text{ Thus} \end{matrix} \quad f^*(y) = 0.$$

Case 3:  $y \geq 1$ . Then the negative of  $y < -1$  is true. Thus

$$f^*(y) = \begin{cases} \infty & y \geq 1 \text{ or } y < -1 \\ 0 & -1 \leq y \leq 1 \end{cases}$$

$$b) f^*(y) = \sup_s (y^T s - \|s\|_1)$$

$$= \sup_s \sum y_i s_i - \sum |s_i|.$$

$$= \sup_s \sum y_i s_i - |s_i|$$

$$= \sum \sup_{s_i} y_i s_i - |s_i|$$

$$= \sum f^*(y_i). \text{ Thus}$$

$$f^*(y) = \begin{cases} \infty & \text{if } \exists y_i < -1 \vee y_i > 1 \\ 0 & \text{if } -1 \leq y_i \leq 1 \end{cases}$$

$$3. p^*(\vec{x}, \vec{c}) = \vec{c}^T \vec{z}_3$$

$$\text{s.t. } \|z_1 - x\|_\infty \leq \epsilon$$

$$(z_{2j}, \hat{z}_{2j}) \in Z_j$$

$$\vec{z}_2 = w_1 \vec{z}_1 + \vec{b}_1$$

$$\hat{z}_3 = w_2 \vec{z}_2 + \vec{b}_2$$

$$I_{B_\epsilon(x)}(z) := \begin{cases} 0 & \text{if } z \in B_\epsilon(x) \\ \infty & \text{otherwise} \end{cases}$$

We can reinforce the constraints by adding  $I_{B_\epsilon(x)}(\vec{z}_1)$  and  $I_{Z_j}(z_{2j}, \hat{z}_{2j})$

since if the constraints are violated then the problem becomes infeasible.

Thus

$$\min_{\vec{z}} \vec{c}^T \vec{z}_3 + I_{B_\epsilon(x)}(\vec{z}_1) + \sum I_{Z_j}(z_{2j}, \hat{z}_{2j})$$

$$\text{s.t. } \vec{z}_2 = w_1 \vec{z}_1 + \vec{b}_1$$

$$\hat{z}_3 = w_2 \vec{z}_2 + \vec{b}_2$$

$$4a) \left\{ \begin{array}{l} C^T \tilde{z}_3 + \mathbb{1}_{B_\epsilon(x)}(\tilde{z}_1) + \sum I_{z_j}(z_{2j}, \tilde{z}_{2j}) \\ \text{s.t. } \tilde{z}_2 = w_1 z_1 + b_1 \\ \tilde{z}_3 = w_2 z_2 + b_2 \end{array} \right\}$$

$$\begin{aligned} &= C^T \tilde{z}_3 + \mathbb{1}_{\frac{\epsilon}{2}(x)}(z_1) + \sum I_{z_j}(z_{2j}, \tilde{z}_{2j}) + V_3^T \tilde{z}_3 - (V_3^T w_2 z_2 - \\ &\quad (V_3^T b_2 + V_2^T \tilde{z}_2) - (C^T w_1 z_1 + V_2^T b_1) \\ &= C^T \tilde{z}_3 + (V_3^T \tilde{z}_3 + \mathbb{1}_{B_\epsilon(x)}(z_1) - V_3^T w_1 z_1 \\ &\quad + \left( \sum_{j=1}^{n_2} \min_{z_{2j}, \tilde{z}_{2j}} (I_z(z_{2j}, \tilde{z}_{2j}) - V_3^T (w_2)_j z_{2j} + V_2^T \tilde{z}_{2j}) \right) \\ &\quad - \sum_{i=1}^2 V_{i+1}^T b_i \end{aligned}$$

4b) Each part of the dual is nicely subdivided according to it's  $z$  or  $\tilde{z}$ . Thus

$$\begin{aligned} g(v_1, v_2) &= \min_{\tilde{z}_3} (C + V_3^T) \tilde{z}_3 + \min_{z_1} \mathbb{1}_{B_\epsilon(x)}(z_1) - (V_2^T w_1 z_1 \\ &\quad + \left( \sum_{j=1}^{n_2} \min_{z_{2j}, \tilde{z}_{2j}} (I_z(z_{2j}, \tilde{z}_{2j}) - V_3^T (w_2)_j z_{2j} + V_2^T \tilde{z}_{2j}) \right) \\ &\quad - \sum_{i=1}^2 V_{i+1}^T b_i \end{aligned}$$

$$q_c) \cdot \min_{\vec{z}_3} (c + v_3)^T \vec{z}_3 = 0 \text{ since } \nabla c = 0 \\ \rightarrow v_3 = -c$$

$$\begin{aligned} & \min_{\vec{z}_1} \left( \mathbb{I}_{B_\epsilon(x)}(\vec{z}_1) - v_3^T w_1 z_1 \right) = -\max_{\vec{z}_1} v_3^T w_1 z_1 - \mathbb{I}_{B_\epsilon(x)}(\vec{z}_1) \\ &= -I_{B_\epsilon}^*(x)(v_2) \\ & \text{Similarly: } \min_{\vec{z}_{2j}, \vec{z}_{2j}} (\mathbb{I}_{\vec{z}}(\vec{z}_{2j}, \vec{z}_{2j}) - v_3^T(w_2)_j z_{2j} + v_{2j} \vec{z}_{2j}) \\ &= -I_{\vec{z}_j}^*(v_{2j}, -v_{2j}) \quad \text{if we define} \\ & \quad v_3^T(w_2)_j = (w_2^T)_j v_3 = \vec{v}_i \end{aligned}$$

Thus L dual is

$$\begin{aligned} & \max_k -I_{B_\epsilon}^*(x)(k_1) + \sum_{j=1}^{n_2} -I_{\vec{z}_j}^*(\vec{v}_{2j}, -v_{2j}) \\ & - \sum_{l=1}^2 v_{l+1}^T b_l \end{aligned}$$

$$\begin{aligned} \text{s.t. } & v_3 = -c \\ & \vec{v}_2 = w_2^T v_3 \\ & \vec{v}_1 = w_1^T v_2 \end{aligned}$$

$$\begin{aligned}
 5. \quad J_{B_\epsilon(x)}^*(\vec{r}) &= \sup_{\substack{s \\ \text{s.t. } \|s-x\|_n \leq \epsilon}} \{ r^T s \} \\
 &= \sup_u r^T x + \epsilon r^T u \quad s = x + \epsilon u \\
 &\quad \text{s.t. } \|u\|_n \leq 1 \\
 &= r^T x + \epsilon \|v\|_1
 \end{aligned}$$

$$\text{Max}_{\boldsymbol{\ell}} \quad -\mathbb{I}_{B_e(x)}^*(\boldsymbol{\ell}_1) + \sum_{j=1}^{n_2} -\mathbb{I}_{Z_j}(\tilde{v}_{2j}, -v_{2j}) \\ - \sum_{i=1}^2 \boldsymbol{\ell}_{i+1}^\top \mathbf{b}_i$$

$$\begin{aligned} \text{s.t. } \quad & \boldsymbol{\nu}_3 = -c \\ & \tilde{\boldsymbol{\nu}}_2 = \mathbf{W}_2^\top \boldsymbol{\nu}_3 \\ & \tilde{\boldsymbol{\nu}}_1 = \mathbf{W}_1^\top \boldsymbol{\nu}_2 \end{aligned}$$

$$\text{Max}_{\boldsymbol{\ell}} \quad -\boldsymbol{\ell}^\top \mathbf{x} - \epsilon \|\boldsymbol{\nu}\|_1 - \sum_{j=1}^{n_2} \boldsymbol{\ell}_{i+1}^\top \mathbf{b}_i + \sum_{j \notin S} l_j \text{ReLU}(\nu_{2j})$$

$$\begin{aligned} \text{s.t. } \quad & \boldsymbol{\nu}_3 = -c \\ & \tilde{\boldsymbol{\nu}}_2 = \mathbf{W}_2^\top \boldsymbol{\nu}_3 \\ & \tilde{\boldsymbol{\nu}}_1 = \mathbf{W}_1^\top \boldsymbol{\nu}_2 \\ & \nu_{2j} = 0 \quad \forall j \in S^- \\ & \nu_{2j} = \tilde{\nu}_{2j} \quad \forall j \in S^+ \\ & \nu_{2j} = \frac{u_j}{l_j - u_j} \quad \forall j \in S \end{aligned}$$

Since we have upper bounded each Feuerle conjugate  $\mathbb{I}_Z^*$  and the bounds are determined by the variables, this is a valid substitution.

$$7. \tilde{z}_2 = w_1 z_1 + b_1$$

We know that  $z_1 = x + \epsilon u$ . Thus

$$\tilde{z}_2 = w_1 x + \epsilon w_1 u + b_1. \text{ Thus}$$

$$\max_{ij} w_{ij}^T u = \|w_1\|_1$$

$$\text{s.t. } \|u\|_\infty \leq 1$$

To maximize/minimize  $\tilde{z}_2$ , you can add/subtract respectively (since the quantity is always positive).

$$\max \tilde{z}_2 = w_1 x + b_1 + \epsilon \|w_1\|_1$$

$$\min \tilde{z}_2 = w_1 x + b_1 - \epsilon \|w_1\|_1$$

8. **Monotonic**: We assume  $y_i \leq \hat{y}_i$  for  $i$  pertaining to the incorrect classes and  $y_{i,\text{true}} \geq \hat{y}_{i,\text{true}}$ . Then

$$\begin{aligned}
 L(y, y_{\text{true}}) &= -\log \left( \frac{e^{\hat{y}_{i,\text{true}}}}{\sum_{j=1}^m e^{\hat{y}_j}} \right) \\
 &= -[\log(e^{\hat{y}_{i,\text{true}}}) - \log(\sum_{j=1}^m e^{\hat{y}_j})] \\
 &\leq \log(\sum_{j=1}^m e^{\hat{y}_j}) - \log(e^{\hat{y}_{i,\text{true}}}) \\
 &\leq \log(\sum_{j=1}^m e^{y_j}) - \log(e^{\hat{y}_{i,\text{true}}}) \\
 &= -\sum_{i \in I} (\hat{y}_{i,\text{true}})_i \log \left( \frac{e^{\hat{y}_i}}{\sum_{j=1}^m e^{\hat{y}_j}} \right)
 \end{aligned}$$

**Translation-Invariant**: We have that

$$\begin{aligned}
 &= -\sum_{i \in I} (\hat{y}_{i,\text{true}})_i \log \left( \frac{e^{\hat{y}_i + \alpha}}{\sum_{j=1}^m e^{\hat{y}_j + \alpha}} \right) \\
 &= -\sum_{i \in I} (\hat{y}_{i,\text{true}})_i \log \left( \frac{e^\alpha \frac{e^{\hat{y}_i}}{\sum_{j=1}^m e^{\hat{y}_j}}}{e^\alpha} \right)
 \end{aligned}$$

$$= - \sum_{i=1}^m (\hat{y}_{true})_i \log \left( \frac{e^{y_i}}{\sum_{j=1}^m e^{y_j}} \right)$$

9. See attached

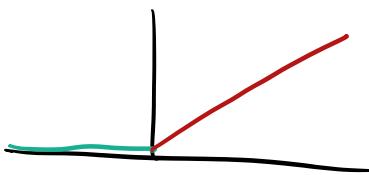
10.

$$f^*(y) = \sup_x \{ y^T x - f(x) \} \quad \forall y \in \mathbb{R}^n$$

What is  $f(x)$ ?

If  $l \leq u \leq 0$ , then the convex hull is

$$\text{thus } f(x) = r > 0$$



$$\sup_{S_1} \{ s_1, \hat{v} \} = 0$$

$$\text{s.t. } s_1 = 0$$

$$(s_1, s_2) \in Z_j$$

Similarly, if  $0 \leq l \leq u$  we have

$$v = \hat{v} \quad \text{and}$$

$$\sup_{S_1} \{ (s_1 - s_2) \hat{v}_j \} = 0 \quad \text{by taking the gradient}$$

$$\text{s.t. } l \leq s_2$$

$$s_2 \leq u$$

$$(s_1, s_2) \in Z_j$$

II. Problem reduces to

$$\sup_s \begin{cases} s_1 \bar{v} - s_2 v \\ \text{s.t.} \\ (s_1, s_2) \in Z_j \end{cases}$$

line: through

$$(l_j, 0); (u_j, u_j)$$

$$\frac{u_j - 0}{u_j - l_j} = m$$

$$y = \frac{u_j}{u_j - l_j}(x - l_j) = \frac{u_j}{u_j - l_j} x - \frac{u_j l_j}{u_j - l_j}$$

$$\rightarrow y = \frac{u_j}{u_j - l_j} x - \frac{u_j l_j}{u_j - l_j}$$

$$s_1 = \frac{u_j s_2}{u_j - l_j} - \frac{u_j l_j}{u_j - l_j}$$

$$\sup_{s_1} \frac{u_j \bar{v} s_2}{u_j - l_j} - \frac{u_j l_j \bar{v}}{u_j - l_j} - s_2 v$$

$$\text{s.t. } (s_1, s_2) \in Z_j$$

Linear, take derivative, we get

$$\frac{u_j \bar{v}}{u_j - l_j} = v \quad \text{and} \quad \sup_{s_1} = -\frac{u_j l_j \bar{v}}{u_j - l_j}$$

$\in -\ell_j v$ . Thus, since the optimum lies at the border, either it is 0 or  $-\ell_j v$ . Thus

$$I_{Z_j}^* = \begin{cases} \text{ReLU}(-\ell_j v) & \text{iff } \frac{u_j v}{u_j - \ell_j} = v \\ +\infty & \text{otherwise} \end{cases}$$