

A. Proofs for Section 3

A.1. Proofs for Section 3.1

A.1.1. PROOF OF THEOREM 1

To start the proof of Theorem 1, we need the following lemma.

Lemma 1. *For any w and r , we have*

$$\sup_{P \in B_{W_\infty}(P_0, r)} R_P(w) = \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r} f(w, \mathbf{x} + \delta) \right]. \quad (16)$$

Proof. Let $T_r^w(\mathbf{x}) = \mathbf{x} + \arg \max_{\{\delta: \|\delta\|_\infty \leq r\}} f(w, \mathbf{x} + \delta)$ with \mathbf{x} is an input data. The existence of $T_r^w(\mathbf{x})$ is guaranteed by the continuity of $f(w, \mathbf{x})$. P_r is the distribution of $T_r^w(\mathbf{x})$ with $\mathbf{x} \sim P_0$. Then

$$\mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r} f(w, \mathbf{x} + \delta) \right] = \mathbb{E}_{P_r} [f(w, \mathbf{x})]. \quad (17)$$

Since

$$W_\infty(P_0, P_r) \leq \mathbb{E}_{P_0} [\|\mathbf{x} - T_r^w(\mathbf{x})\|_\infty] \leq r, \quad (18)$$

we have

$$\mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r} f(w, \mathbf{x} + \delta) \right] \leq \sup_{P \in B_{W_\infty}(P_0, r)} R_P(w). \quad (19)$$

On the other hand, let $P_r^* \in \arg \max_{P \in B_{W_\infty}(P_0, r)} R_P(w)$. Due to Kolmogorov's theorem, P_r^* can be distribution of some random vector \mathbf{z} , due to the definition of W_∞ -distance, we have $\|\mathbf{z} - \mathbf{x}\|_\infty \leq r$ holds almost surely. Then we conclude

$$\sup_{P \in B_{W_\infty}(P_0, r)} R_P(w) = R_{P_r^*}(w) = \mathbb{E}_{P_r^*} [f(w, \mathbf{z})] \leq \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r} f(w, \mathbf{x} + \delta) \right]. \quad (20)$$

Thus, we get the conclusion. \square

This lemma shows that the distributional perturbation measured by W_∞ -distance is equivalent to input perturbation. Hence we can study W_{inf} -distributional robustness through ℓ_{inf} -input-robustness. The basic tool for our proof is the covering number, which is defined as follows.

Definition 2. (Wainwright, 2019) *A r -cover of $(\mathcal{X}, \|\cdot\|_p)$ is any point set $\{\mathbf{u}_i\} \subseteq \mathcal{X}$ such that for any $\mathbf{u} \in \mathcal{X}$, there exists \mathbf{u}_i satisfies $\|\mathbf{u} - \mathbf{u}_i\|_p \leq r$. The covering number $\mathcal{N}(r, \mathcal{X}, \|\cdot\|_p)$ is the cardinality of the smallest r -cover.*

Now we are ready to give the proof of Theorem 1 which is motivated by (Xu & Mannor, 2012).

Proof of Theorem 1. We can construct a r -cover to $(\mathcal{X}, \|\cdot\|_2)$ then $\mathcal{N}(r, \mathcal{X}, \|\cdot\|_2) \leq (2d_0)^{(2D/r^2+1)} = N$, because the \mathcal{X} can be covered by a polytope with ℓ_2 -diameter smaller than $2D$ and $2d_0$ vertices, see (Vershynin, 2018) Theorem 0.0.4 for details. Due to the geometrical structure, we have $\mathcal{N}(r, \mathcal{X}, \|\cdot\|_\infty) \leq (2d_0)^{(2D/r^2+1)}$. Then, there exists (C_1, \dots, C_N) covers $(\mathcal{X}, \|\cdot\|_\infty)$ where C_i is disjoint with each other, and $\|\mathbf{u} - \mathbf{v}\|_\infty \leq r$ for any $\mathbf{u}, \mathbf{v} \in C_i$. This can be constructed by $C_i = \hat{C}_i \cap \left(\bigcup_{j=1}^{i-1} \hat{C}_j \right)^c$ with $(\hat{C}_1, \dots, \hat{C}_N)$ covers $(\mathcal{X}, \|\cdot\|_\infty)$, and the diameter of each \hat{C}_i is smaller than r since

$\mathcal{N}(r, \mathcal{X}, \|\cdot\|_\infty) \leq N$. Let $A_j = \{\mathbf{x}_i : \mathbf{x}_i \in C_j\}$, and $|A_j|$ is the cardinality of A_j . Due to Lemma 1, we have

$$\begin{aligned}
 \left| \sup_{P \in B_{W_\infty}(P_0, r_0)} R_P(\mathbf{w}) - R_{P_n}(\mathbf{w}) \right| &= \left| \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r_0} f(\mathbf{w}, \mathbf{x} + \delta) \right] - R_{P_n}(\mathbf{w}) \right| \\
 &= \left| \sum_{j=1}^N \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r_0} f(\mathbf{w}, \mathbf{x} + \delta) \mid \mathbf{x} \in C_j \right] P_0(C_j) - R_{P_n}(\mathbf{w}) \right| \\
 &\leq \left| \sum_{j=1}^N \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r_0} f(\mathbf{w}, \mathbf{x} + \delta) \mid \mathbf{x} \in C_j \right] \frac{|A_j|}{n} - \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}, \mathbf{x}_i) \right| \\
 &\quad + \left| \sum_{j=1}^N \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r_0} f(\mathbf{w}, \mathbf{x} + \delta) \mid \mathbf{x} \in C_j \right] \left(\frac{|A_j|}{n} - P_0(C_j) \right) \right| \tag{21} \\
 &\leq \left| \frac{1}{n} \sum_{j=1}^N \sum_{\mathbf{x}_i \in C_j} \sup_{\mathbf{x} \in C_j + B_\infty(\mathbf{0}, r_0)} |f(\mathbf{w}, \mathbf{x}) - f(\mathbf{w}, \mathbf{x}_i)| \right| + M \sum_{j=1}^N \left| \frac{|A_j|}{n} - P_0(C_j) \right| \\
 &\stackrel{a}{\leq} \frac{1}{n} \sum_{i=1}^n \sup_{\|\delta\|_\infty \leq 2r} |f(\mathbf{w}, \mathbf{x}_i + \delta) - f(\mathbf{w}, \mathbf{x}_i)| + M \sum_{j=1}^N \left| \frac{|A_j|}{n} - P_0(C_j) \right| \\
 &\leq \epsilon + M \sum_{j=1}^N \left| \frac{|A_j|}{n} - P_0(C_j) \right|.
 \end{aligned}$$

Here a is due to $C_j + B_\infty(\mathbf{0}, r) \subseteq B_\infty(\mathbf{x}_i, 2r)$ when $\mathbf{x}_i \in C_j$, since ℓ_∞ -diameter of C_j is smaller than r . The last inequality is due to $(2r, \epsilon, P_n, \infty)$ -robustness of $f(\mathbf{w}, \mathbf{x})$. On the other hand, due to Proposition A6.6 in (van der Vaart & Wellner, 2000), we have

$$\mathbb{P} \left(\sum_{j=1}^N \left| \frac{|A_j|}{n} - P_0(C_j) \right| \geq \theta \right) \leq 2^N \exp \left(\frac{-n\theta^2}{2} \right). \tag{22}$$

Combine this with (21), due to the value of N , we get the conclusion. \square

A.1.2. PROOF OF THEOREM 2

There is a little difference of proving Theorem 2 compared with Theorem 1. Because the out-distribution P constrained in $B_{W_\infty}(P_0, r)$ only correspond with OOD data that contained in a ℓ_∞ -ball of in-distribution data almost surely, see Lemma 1 for a rigorous description. Hence, we can utilize ℓ_∞ -robustness of model to derive the OOD generalization under W_∞ -distance by Theorem 1. However, in the regime of W_2 -distance, roughly speaking, the transformed OOD data $T_r^w(\mathbf{x})$ is contained in a ℓ_2 -ball of \mathbf{x} in expectation. Thus, Lemma 1 is invalid under W_2 -distance.

To discuss the OOD generalization under W_2 -distance, we need to give a delicate characterization to the distribution $P \in B_{W_2}(P_0, r)$. First, we need the following lemma.

Lemma 2. *For any r and \mathbf{w} , let $P_r^* \in \arg \max_{P \in B_{W_2}(P_0, r)} R_P(\mathbf{w})$. Then, there exists a mapping $T_r^w(\mathbf{x})$ such that $T_r^w(\mathbf{x}) \sim P_r^*$ with $\mathbf{x} \sim P_0$.*

Proof. The proof of Theorem 6 in (Sinha et al., 2018) shows that

$$R_{P_r^*}(\mathbf{w}) = \sup_{P \in B_{W_2}(P_0, r)} R_P(\mathbf{w}) = \inf_{\lambda \geq 0} \sup_{P, \pi \in (P, P_0)} \left(\int_{\mathcal{X} \times \mathcal{X}} f(\mathbf{w}, \mathbf{x}) - \lambda \|\mathbf{x} - \mathbf{z}\|^2 d\pi(\mathbf{x}, \mathbf{z}) + \lambda r \right). \tag{23}$$

We next show that the supremum over π in the last equality is attained by the joint distribution $(T_r^w(\mathbf{x}), \mathbf{x})$, which implies our conclusion. For any $\lambda > 0$, we have

$$\sup_{P, \pi \in (P, P_0)} \left(\int_{\mathcal{X} \times \mathcal{X}} f(\mathbf{w}, \mathbf{x}) - \lambda \|\mathbf{x} - \mathbf{z}\|^2 d\pi(\mathbf{x}, \mathbf{z}) \right) \leq \int_{\mathcal{X}} \sup_{\mathbf{x}} (f(\mathbf{w}, \mathbf{x}) - \lambda \|\mathbf{x} - \mathbf{z}\|^2) dP_0(\mathbf{z}), \tag{24}$$

due to the supremum in the left hand side is taken over P and π . On the other hand, let $P(\cdot | \mathbf{z})$ and $\mathbf{x}(\cdot)$ respectively be the regular conditional distribution on \mathcal{X} with \mathbf{z} given and the function on \mathcal{X} . Since $P(\cdot | \mathbf{z})$ is measurable,

$$\begin{aligned} \sup_{P, \pi \in (P, P_0)} \left(\int_{\mathcal{X} \times \mathcal{X}} f(\mathbf{w}, \mathbf{x}) - \lambda \|\mathbf{x} - \mathbf{z}\|^2 d\pi(\mathbf{x}, \mathbf{z}) \right) &\geq \sup_{P(\cdot | \mathbf{z})} \left(\int_{\mathcal{X} \times \mathcal{X}} f(\mathbf{w}, \mathbf{x}) - \lambda \|\mathbf{x} - \mathbf{z}\|^2 dP(\mathbf{x} | \mathbf{z}) dP_0(\mathbf{z}) \right) \\ &\geq \sup_{\mathbf{x}(\cdot)} \left(\int_{\mathcal{X}} f(\mathbf{w}, \mathbf{x}(\mathbf{z})) - \lambda \|\mathbf{x}(\mathbf{z}) - \mathbf{z}\|^2 dP_0(\mathbf{z}) \right) \\ &\geq \int_{\mathcal{X}} \sup_{\mathbf{x}} (f(\mathbf{w}, \mathbf{x}) - \lambda \|\mathbf{x} - \mathbf{z}\|^2) dP_0(\mathbf{z}). \end{aligned} \quad (25)$$

Thus, we get the conclusion. \square

Proof of Theorem 2. Similar to the proof of Theorem 1, we can construct a disjoint cover (C_1, \dots, C_N) to $(\mathcal{X}, \|\cdot\|_2)$ such that $N \leq (2d_0)^{(2\epsilon^2 D/r^2 + 1)}$, and the l_2 -diameter of each C_i is smaller than r/ϵ . Let $P_r^* \in \arg \max_{P \in B_{W_2}(P_0, r)} R_P(\mathbf{w})$, by Lemma 2, we have

$$\begin{aligned} \sup_{P \in B_{W_2}(P_0, r)} R_P(\mathbf{w}) &= R_{P_r^*}(\mathbf{w}) \\ &= \mathbb{E}_{P_0} [f(\mathbf{w}, T_r^{\mathbf{w}}(\mathbf{x}))] \\ &= \mathbb{E}_{P_0} [f(\mathbf{w}, T_r^{\mathbf{w}}(\mathbf{x})) (\mathbf{1}_{T_r^{\mathbf{w}}(\mathbf{x}) \in B_2(\mathbf{x}, r/\epsilon)} + \mathbf{1}_{T_r^{\mathbf{w}}(\mathbf{x}) \notin B_2(\mathbf{x}, r/\epsilon)})] \\ &\leq \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_2 \leq r/\epsilon} f(\mathbf{w}, \mathbf{x} + \delta) \right] + M \mathbb{P}(T_r^{\mathbf{w}}(\mathbf{x}) \notin B_2(\mathbf{x}, r/\epsilon)). \end{aligned} \quad (26)$$

Due to the definition of $T_r^{\mathbf{w}}(\mathbf{x})$, by Markov's inequality, we have

$$\left(\frac{r}{\epsilon} \right) \mathbb{P}(T_r^{\mathbf{w}}(\mathbf{x}) \notin B_2(\mathbf{x}, r/\epsilon)) \leq \int_{\mathcal{X}} \|T_r^{\mathbf{w}}(\mathbf{x}) - \mathbf{x}\|^2 dP_0(\mathbf{x}) = W_2(P_0, P_r^*) \leq r. \quad (27)$$

Plugging this into (26), and due to the definition of Wasserstein distance, we have

$$\mathbb{E}_{P_0} \left[\sup_{\|\delta\|_2 \leq r/\epsilon} f(\mathbf{w}, \mathbf{x} + \delta) \right] \leq \sup_{P \in B_{W_2}(P_0, r)} R_P(\mathbf{w}) \leq \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_2 \leq r/\epsilon} f(\mathbf{w}, \mathbf{x} + \delta) \right] + M\epsilon. \quad (28)$$

Similar to the proof of Theorem 1, due to the model is $(2r/\epsilon, \epsilon, P_n, 2)$ -robust, we have

$$\left| \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_2 \leq r/\epsilon} f(\mathbf{w}, \mathbf{x} + \delta) \right] - R_{P_n}(\mathbf{w}) \right| \leq \epsilon + M \sqrt{\frac{(2d_0)^{(2\epsilon^2 D/r^2 + 1)} \log 2 + 2 \log(1/\theta)}{n}} \quad (29)$$

holds with probability at least $1 - \theta$. Combining this with (28), we get the conclusion. \square

A.2. Proofs for Section 3.2

The proof of Theorem 3 is same for $p \in \{2, \infty\}$, we take $p = \infty$ as an example. Before providing the proof, we first give a lemma to characterize the convergence rate of the first inner loop in Algorithm 1.

Lemma 3. For any $\mathbf{w}, \mathbf{x} \in \{\mathbf{x}_i\}$, and r , there exists $\delta^* \in \arg \max_{\{\delta: \|\delta\|_\infty \leq r\}} f(\mathbf{w}, \mathbf{x} + \delta)$ such that

$$\|\delta_{K+1} - \delta^*\|^2 \leq \left(1 - \frac{\mu_{\mathbf{x}}}{L_{22}} \right)^K \|\delta_1 - \delta^*\|^2 \quad (30)$$

when $\delta_{k+1} = \text{Proj}_{B_\infty(\mathbf{0}, r)}(\delta_k + \eta_{\mathbf{x}} \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x} + \delta_k))$ with $\eta_{\mathbf{x}} = 1/L_{22}$.

Proof. The existence of δ^* is due to the continuity of $f(\mathbf{w}, \cdot)$. Then

$$\begin{aligned}
 f(\mathbf{w}, \mathbf{x} + \delta^*) - f(\mathbf{w}, \mathbf{x} + \delta_{k+1}) &= f(\mathbf{w}, \mathbf{x} + \delta^*) - f(\mathbf{w}, \mathbf{x} + \delta_k) + f(\mathbf{w}, \mathbf{x} + \delta_k) - f(\mathbf{w}, \mathbf{x} + \delta_{k+1}) \\
 &\stackrel{a}{\leq} \langle \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x} + \delta_k), \delta^* - \delta_k \rangle - \frac{\mu_{\mathbf{x}}}{2} \|\delta_k - \delta^*\|^2 \\
 &\quad + \langle \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x} + \delta_k), \delta_k - \delta_{k+1} \rangle + \frac{L_{22}}{2} \|\delta_{k+1} - \delta_k\|^2 \\
 &= \langle \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x} + \delta_k), \delta^* - \delta_{k+1} \rangle - \frac{\mu_{\mathbf{x}}}{2} \|\delta_k - \delta^*\|^2 + \frac{L_{22}}{2} \|\delta_{k+1} - \delta_k\|^2 \\
 &\stackrel{b}{\leq} L_{22} \langle \delta_{k+1} - \delta_k, \delta^* - \delta_{k+1} \rangle - \frac{\mu_{\mathbf{x}}}{2} \|\delta_k - \delta^*\|^2 + \frac{L_{22}}{2} \|\delta_{k+1} - \delta_k\|^2 \\
 &= L_{22} \langle \delta_{k+1} - \delta_k, \delta^* - \delta_k \rangle - \frac{\mu_{\mathbf{x}}}{2} \|\delta_k - \delta^*\|^2 - \frac{L_{22}}{2} \|\delta_{k+1} - \delta_k\|^2,
 \end{aligned} \tag{31}$$

where a is due to the L_{22} -Lipschitz continuity of $\nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x})$ and strongly convexity, b is because the property of projection (see Lemma 3.1 in (Bubeck, 2014)). Then we get

$$\begin{aligned}
 \|\delta_{k+1} - \delta^*\|^2 &= \|\delta_{k+1} - \delta_k\|^2 + \|\delta_k - \delta^*\|^2 + 2\langle \delta_{k+1} - \delta_k, \delta_k - \delta^* \rangle \\
 &\leq \left(1 - \frac{\mu_{\mathbf{x}}}{L_{22}}\right) \|\delta_k - \delta^*\|^2
 \end{aligned} \tag{32}$$

by plugging (31) into the above equality and $f(\mathbf{w}, \mathbf{x} + \delta^*) - f(\mathbf{w}, \mathbf{x} + \delta_{k+1}) \geq 0$. Thus, we get the conclusion. \square

This lemma shows that the inner loop in Algorithm 1 can efficiently approximate the worst-case perturbation for any \mathbf{w}_t and \mathbf{x}_i . Now we are ready to give the proof of Theorem 3.

We need the following lemma, which is Theorem 6 in (Rakhlin et al., 2012).

Lemma 4. *Let $\{\xi_1, \dots, \xi_t\}$ be a martingale difference sequence with a uniform upper bound b . Let $V_t = \sum_{j=1}^t \text{Var}(\xi_j | \mathcal{F}_{j-1})$ with \mathcal{F}_j is the σ -field generated by $\{\xi_1, \dots, \xi_j\}$. Then for every a and $v > 0$,*

$$\mathbb{P} \left(\bigcup_{s \leq t} \left(\left\{ \sum_{j=1}^s \xi_j \geq a \right\} \cap \{V_t \leq v\} \right) \right) \leq \exp \left(\frac{-a^2}{2(v + ba)} \right). \tag{33}$$

This is a type of Bennett's inequality which is sharper compared with Azuma-Hoeffding's inequality when the variance v is much smaller than uniform bound b .

A.2.1. PROOF OF THEOREM 3

Proof. With a little abuse of notation, let $r(p) = r$ and define $g(\mathbf{w}, \mathbf{x}) = \sup_{\delta: \|\delta\|_{\infty} \leq r} f(\mathbf{w}, \mathbf{x} + \delta)$. Lemma A.5 in (Nouiehed et al., 2019) implies $g(\mathbf{w}, \mathbf{x})$ has $L_{11} + \frac{L_{12}L_{21}}{\mu_{\mathbf{x}}}$ -Lipschitz continuous gradient with respect to \mathbf{w} for any specific \mathbf{x} . Then $\tilde{R}_{P_n}(\mathbf{w})$ has $L = L_{11} + \frac{L_{12}L_{21}}{\mu_{\mathbf{x}}}$ -Lipschitz continuous gradient. Let $\mathbf{x}^* \in \mathbf{x} + \arg \max_{\delta: \|\delta\|_{\infty} \leq r} f(\mathbf{w}, \mathbf{x} + \delta)$, due to the Lipschitz gradient of $\tilde{R}_{P_n}(\mathbf{w})$,

$$\begin{aligned}
 \tilde{R}_{P_n}(\mathbf{w}_{t+1}) - \tilde{R}_{P_n}(\mathbf{w}_t) &\leq \langle \nabla \tilde{R}_{P_n}(\mathbf{w}_t), \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{L}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\
 &= -\eta_{\mathbf{w}_t} \langle \nabla \tilde{R}_{P_n}(\mathbf{w}_t), \nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_{i_t} + \delta_K) \rangle + \frac{\eta_{\mathbf{w}_t}^2 L}{2} \|\nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_{i_t} + \delta_K)\|^2 \\
 &= -\eta_{\mathbf{w}_t} \|\nabla \tilde{R}_{P_n}(\mathbf{w}_t)\|^2 + \eta_{\mathbf{w}_t} \langle \nabla \tilde{R}_{P_n}(\mathbf{w}_t), \nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_{i_t}^*) - \nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_{i_t} + \delta_K) \rangle \\
 &\quad + \eta_{\mathbf{w}_t} \langle \nabla \tilde{R}_{P_n}(\mathbf{w}_t), \nabla \tilde{R}_{P_n}(\mathbf{w}_t) - \nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_{i_t}^*) \rangle + \frac{\eta_{\mathbf{w}_t}^2 L}{2} \|\nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_{i_t} + \delta_K)\|^2.
 \end{aligned} \tag{34}$$

Here the last equality is due to $\nabla_{\mathbf{w}} g(\mathbf{w}, \mathbf{x}) = \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}^*)$ (Similar to Danskin's theorem, see Lemma A.5 in (Nouiehed et al., 2019)), and $\mathbf{x}_{i_t}^*$ is the local maxima approximated by $\mathbf{x}_{i_t} + \delta_K$ in Lemma 3. By taking expectation to \mathbf{w}_{t+1} with \mathbf{w}_t

given in the both side of the above equation, Jensen's inequality, combining Lemma 3 and $\eta_{w_t} = 1/\mu_w t$,

$$\begin{aligned}
 \mathbb{E}[\tilde{R}_{P_n}(w_{t+1})] - \tilde{R}_{P_n}(w^*) &\leq \tilde{R}_{P_n}(w_t) - \tilde{R}_{P_n}(w^*) - \eta_{w_t} \|\nabla \tilde{R}_{P_n}(w_t)\|^2 \\
 &\quad + \mathbb{E} \left[\eta_{w_t} \|\nabla \tilde{R}_{P_n}(w_t)\| \|\nabla_w f(w_t, x_{i_t}^*) - \nabla_w f(w, x_{i_t} + \delta_K)\| \right] + \frac{\eta_{w_t}^2 G^2 L}{2} \\
 &\leq \tilde{R}_{P_n}(w_t) - \tilde{R}_{P_n}(w^*) - \eta_{w_t} \|\nabla \tilde{R}_{P_n}(w_t)\|^2 + \eta_{w_t} \|\nabla \tilde{R}_{P_n}(w_t)\| \left(1 - \frac{\mu_x}{L_{22}}\right)^K \mathbb{E} [\|\delta_1 - \delta_{i_t}^*\|^2] + \frac{\eta_{w_t}^2 G^2 L}{2} \\
 &\leq (1 - 2\mu_w \eta_{w_t}) \left(\tilde{R}_{P_n}(w_t) - \tilde{R}_{P_n}(w^*) \right) + \eta_{w_t}^2 G^2 L \\
 &= \left(1 - \frac{2}{t}\right) \left(\tilde{R}_{P_n}(w_t) - \tilde{R}_{P_n}(w^*) \right) + \frac{G^2 L}{\mu_w^2 t^2}.
 \end{aligned} \tag{35}$$

Here the third inequality is because

$$\eta_{w_t} \|\nabla \tilde{R}_{P_n}(w_t)\| \left(1 - \frac{\mu_x}{L_{22}}\right)^K \|\delta_1 - \delta_{i_t}^*\|^2 \leq \eta_{w_t} G \left(1 - \frac{\mu_x}{L_{22}}\right)^K 4d_0 r^2 \leq \frac{\eta_{w_t}^2 G^2 L}{2}, \tag{36}$$

for any $\delta_{i_t}^*$, since

$$K \log \left(1 - \frac{\mu_x}{L_{22}}\right) \leq -K \frac{\mu_x}{L_{22}} \leq \log \left(\frac{GL}{8T\mu_w d_0 r^2}\right). \tag{37}$$

Then by induction,

$$\begin{aligned}
 \mathbb{E}[\tilde{R}_{P_n}(w_{t+1})] - \tilde{R}_{P_n}(w^*) &\leq \frac{G^2 L}{\mu_w^2} \sum_{j=2}^t \frac{1}{j^2} \prod_{k=j+1}^t \left(1 - \frac{2}{k}\right) \\
 &= \frac{G^2 L}{\mu_w^2} \sum_{j=2}^t \frac{1}{j^2} \frac{(j-1)j}{(t-1)t} \\
 &\leq \frac{G^2 L}{t\mu_w^2}.
 \end{aligned} \tag{38}$$

Thus we get the first conclusion of convergence in expectation by taking $t = T$ for $t \geq 2$. For the second conclusion, let us define $\xi_t = \langle \nabla \tilde{R}_{P_n}(w_t), \nabla \tilde{R}_{P_n}(w_t) - \nabla_w f(w_t, x_{i_t}^*) \rangle$. Then Schwarz inequality implies that

$$|\xi_t| \leq \|\nabla \tilde{R}_{P_n}(w_t)\| \|\nabla \tilde{R}_{P_n}(w_t) - \nabla_w f(w_t, x_{i_t}^*)\| \leq 2G^2. \tag{39}$$

Similar to (35), for $t \geq 2$,

$$\begin{aligned}
 \tilde{R}_{P_n}(w_{t+1}) - \tilde{R}_{P_n}(w^*) &\leq (1 - 2\mu_w \eta_{w_t}) \left(\tilde{R}_{P_n}(w_t) - \tilde{R}_{P_n}(w^*) \right) + \eta_{w_t}^2 G^2 L + 2\eta_{w_t} \xi_t \\
 &\leq \frac{G^2 L}{t\mu_w^2} + \frac{2}{\mu_w} \sum_{j=2}^t \frac{\xi_j}{j} \prod_{k=j+1}^t \left(1 - \frac{2}{k}\right) \\
 &= \frac{G^2 L}{t\mu_w^2} + \frac{2}{\mu_w} \sum_{j=2}^t \frac{1}{j} \frac{(j-1)j}{(t-1)t} \xi_j \\
 &= \frac{G^2 L}{t\mu_w^2} + \frac{2}{\mu_w} \sum_{j=2}^t \frac{(j-1)}{(t-1)t} \xi_j.
 \end{aligned} \tag{40}$$

Since the second term in the last inequality is upper bonded by $\sum_{j=2}^t \xi_j$ which is a sum of martingale difference, and $|\xi_j| \leq 2G^2$, a simple Azuma-Hoeffding's inequality based on bounded martingale difference (Corollary 2.20 in (Wainwright, 2019)) can give a $\mathcal{O}(1/\sqrt{t})$ convergence rate in the high probability. However, we can sharpen the convergence rate via a Bennett's inequality (Proposition 3.19 in (Duchi, 2016)), because the conditional variance of ξ_j will decrease across training. We consider the conditional variance of $\sum_{j=2}^t (j-1)\xi_j$, let \mathcal{F}_j be the σ -field generated by $\{w_1, \dots, w_j\}$, since $\mathbb{E}[\xi_j] = 0$

we have

$$\begin{aligned}
 \text{Var} \left(\sum_{j=2}^t (j-1) \xi_j \mid \mathcal{F}_{j-1} \right) &= \sum_{j=2}^t (j-1)^2 \text{Var} (\xi_j \mid \mathcal{F}_{j-1}) \\
 &= \sum_{j=2}^t (j-1)^2 \mathbb{E} [\xi_j^2 \mid \mathcal{F}_{j-1}] \\
 &\leq 4G^2 \sum_{j=2}^t (j-1)^2 \|\nabla \tilde{R}_{P_n}(\mathbf{w}_j)\|^2 \\
 &\leq 8G^2 L \sum_{j=2}^t (j-1)^2 \left(\tilde{R}_{P_n}(\mathbf{w}_j) - \tilde{R}_{P_n}(\mathbf{w}^*) \right),
 \end{aligned} \tag{41}$$

where first inequality is from Schwarz's inequality and the last inequality is because

$$\begin{aligned}
 \tilde{R}_{P_n}(\mathbf{w}^*) - \tilde{R}_{P_n}(\mathbf{w}) &\leq \tilde{R}_{P_n} \left(\mathbf{w} - \frac{1}{L} \nabla \tilde{R}_{P_n}(\mathbf{w}) \right) - \tilde{R}_{P_n}(\mathbf{w}) \\
 &\leq - \left\langle \nabla \tilde{R}_{P_n}(\mathbf{w}), \frac{1}{L} \nabla \tilde{R}_{P_n}(\mathbf{w}) \right\rangle + \frac{L}{2} \left\| \frac{1}{L} \nabla \tilde{R}_{P_n}(\mathbf{w}) \right\|^2 \\
 &= - \frac{1}{2L} \left\| \nabla \tilde{R}_{P_n}(\mathbf{w}) \right\|^2,
 \end{aligned} \tag{42}$$

for any \mathbf{w} . By applying Lemma 4, as long as $T \geq 4$ and $0 < \theta < 1/e$, then with probability at least $1 - \theta$, for all $t \leq T$,

$$\begin{aligned}
 &\tilde{R}_{P_n}(\mathbf{w}_{t+1}) - \tilde{R}_{P_n}(\mathbf{w}^*) \\
 &\leq \frac{8G}{\mu_{\mathbf{w}}(t-1)t} \max \left\{ \sqrt{2L \sum_{j=2}^t (j-1)^2 \left(\tilde{R}_{P_n}(\mathbf{w}_j) - \tilde{R}_{P_n}(\mathbf{w}^*) \right)}, G(t-1) \sqrt{\log \left(\frac{\log T}{\theta} \right)} \right\} \sqrt{\log \left(\frac{\log T}{\theta} \right)} + \frac{G^2 L}{t\mu_{\mathbf{w}}^2} \\
 &\leq \frac{8G \sqrt{\log(\log(T/\theta))}}{\mu_{\mathbf{w}}(t-1)t} \sqrt{2L \sum_{j=2}^t (j-1)^2 \left(\tilde{R}_{P_n}(\mathbf{w}_j) - \tilde{R}_{P_n}(\mathbf{w}^*) \right)} + \frac{(8\mu_{\mathbf{w}} G^2 \log(\log(T/\theta)) + G^2 L)}{t\mu_{\mathbf{w}}^2}.
 \end{aligned} \tag{43}$$

Then, an upper bound to the first term in the last inequality can give our conclusion. Note that if $\tilde{R}_{P_n}(\mathbf{w}_j) - \tilde{R}_{P_n}(\mathbf{w}^*)$ is smaller than $\mathcal{O}(1/j-1)$, the conclusion is full-filled. To see this, we should find a large constant a such that $\tilde{R}_{P_n}(\mathbf{w}_{t+1}) - \tilde{R}_{P_n}(\mathbf{w}^*) \leq a/t$. This is clearly hold when $a \geq G^2/2\mu_{\mathbf{w}}$ for $t = 1$ due to the PL inequality and bounded gradient. For $t \geq 2$, we find this a by induction. Let $b = 8G\sqrt{2L \log(\log(T/\theta))}/\mu_{\mathbf{w}}$ and $c = (8\mu_{\mathbf{w}} G^2 \log(\log(T/\theta)) + G^2 L)/\mu_{\mathbf{w}}^2$. A satisfactory a yields

$$\frac{a}{t} \geq \frac{b}{(t-1)t} \sqrt{a \sum_{j=2}^t (j-1)} + \frac{c}{t} = \frac{b}{(t-1)t} \sqrt{\frac{at(t-1)}{2}} + \frac{c}{t} \geq \frac{1}{t} \left(b\sqrt{\frac{a}{2}} + c \right). \tag{44}$$

By solving a quadratic inequality, we conclude that $a - b\sqrt{a/2} - c \geq 0$. Then

$$a \geq \left(\frac{b + \sqrt{b^2 + 8c}}{2\sqrt{2}} \right)^2. \tag{45}$$

By taking

$$a \geq 2 \left(\frac{2b^2 + 8c}{8} \right) \geq \left(\frac{b + \sqrt{b^2 + 8c}}{2\sqrt{2}} \right)^2, \tag{46}$$

we get

$$a \geq \frac{64G^2 L \log(\log(T/\theta))}{\mu_{\mathbf{w}}^2} + \frac{(16\mu_{\mathbf{w}} G^2 \log(\log(T/\theta)) + G^2 L)}{\mu_{\mathbf{w}}^2} = \frac{G^2 \log(\log(T/\theta))(64L + 16\mu_{\mathbf{w}}) + G^2 L}{\mu_{\mathbf{w}}^2}, \tag{47}$$

due to the value of b and c . Hence, we get the conclusion by taking $t = T$. \square

A.2.2. PROOF OF PROPOSITION 1

Proof. From the definition of $\tilde{R}_{P_n}(\mathbf{w})$, for any $r \geq 0$, we have

$$\frac{1}{n} \sum_{i=1}^n \sup_{\|\delta\|_p \leq r} (f(\mathbf{w}, \mathbf{x}_i + \delta) - f(\mathbf{w}, \mathbf{x}_i)) \leq \tilde{R}_{P_n}(\mathbf{w}) \leq \epsilon. \quad (48)$$

On the other hand

$$\frac{1}{n} \sum_{i=1}^n \sup_{\|\delta\|_p \leq r} (f(\mathbf{w}, \mathbf{x}_i) - f(\mathbf{w}, \mathbf{x}_i + \delta)) \leq R_{P_n}(\mathbf{w}) \leq \tilde{R}_{P_n}(\mathbf{w}) \leq \epsilon. \quad (49)$$

Take a sum to the two above inequalities, we get

$$\frac{1}{n} \sum_{i=1}^n \sup_{\|\delta\|_p \leq r} |f(\mathbf{w}, \mathbf{x}_i + \delta) - f(\mathbf{w}, \mathbf{x}_i)| \leq \frac{1}{n} \sum_{i=1}^n \left(\sup_{\|\delta\|_p \leq r} f(\mathbf{w}, \mathbf{x}_i + \delta) - \inf_{\|\delta\|_p \leq r} f(\mathbf{w}, \mathbf{x}_i + \delta) \right) \leq 2\epsilon. \quad (50)$$

Then the conclusion is verified. \square

B. Proofs for Section 4

B.1. Proof of Theorem 4

Proof. We have $r(\infty) = r$ in this theorem. The key is to bound the $|\sup_{P \in B_{W_\infty}(P_0, r)} R_P(\mathbf{w}_{\text{pre}}) - \sup_{Q \in B_{W_\infty}(Q_0, r)} R_Q(\mathbf{w}_{\text{pre}})|$, then triangle inequality and Hoeffding's inequality imply the conclusion. Let $P_r^* \in \arg \max_{P \in B_{W_\infty}(P_0, r)} R_P(\mathbf{w}_{\text{pre}})$. For any given \mathbf{x} , due to the continuity of $f(\mathbf{w}_{\text{pre}}, \cdot)$, similar to Lemma 1, we can find the $T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x}) = \mathbf{x} + \arg \max_{\{\delta: \|\delta\|_\infty \leq r\}} f(\mathbf{w}_{\text{pre}}, \mathbf{x} + \delta)$. Then due to Lemma 1,

$$R_{P_r^*}(\mathbf{w}_{\text{pre}}) = \mathbb{E}_{P_0} \left[\sup_{\|\delta\|_\infty \leq r} f(\mathbf{w}_{\text{pre}}, \mathbf{x} + \delta) \right]. \quad (51)$$

Thus, $T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x}) \sim P_r^*$ when $\mathbf{x} \sim P_0$. We can find $\mathbf{z} \sim Q_0$ due to the Kolmogorov's Theorem, and let $T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{z}) \sim Q_r^*$. By the definition of W_∞ -distance, one can verify $W_\infty(Q_0, Q_r^*) \leq r$ as well as $R_{Q_r^*}(\mathbf{w}_{\text{pre}}) \leq \epsilon_{\text{pre}}$. Note that $0 \leq f(\mathbf{w}_{\text{pre}}, \cdot) \leq M$, then

$$\begin{aligned} |R_{P_r^*}(\mathbf{w}_{\text{pre}}) - R_{Q_r^*}(\mathbf{w}_{\text{pre}})| &= \left| \int_{\mathcal{X}} f(\mathbf{w}_{\text{pre}}, \mathbf{x}) dP_r^*(\mathbf{x}) - \int_{\mathcal{X}} f(\mathbf{w}_{\text{pre}}, \mathbf{x}) dQ_r^*(\mathbf{x}) \right| \\ &= \left| \int_{\mathcal{X}} f(\mathbf{w}_{\text{pre}}, T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x})) dP_0(\mathbf{x}) - \int_{\mathcal{X}} f(\mathbf{w}_{\text{pre}}, T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x})) dQ_0(\mathbf{x}) \right| \\ &\leq \int_{\mathcal{X}} |f(\mathbf{w}_{\text{pre}}, T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x}))| |dP_0(\mathbf{x}) - dQ_0(\mathbf{x})| \\ &\leq M \int_{\mathcal{X}} |dP_0(\mathbf{x}) - dQ_0(\mathbf{x})| \\ &= 2MTV(P_0, Q_0). \end{aligned} \quad (52)$$

The last equality is from the definition of total variation distance (Villani, 2008). Thus a simple triangle inequality implies that

$$R_{P_r^*}(\mathbf{w}_{\text{pre}}) \leq |R_{P_r^*}(\mathbf{w}_{\text{pre}}) - R_{Q_r^*}(\mathbf{w}_{\text{pre}})| + R_{Q_r^*}(\mathbf{w}_{\text{pre}}) \leq \epsilon_{\text{pre}} + 2MTV(P_0, Q_0). \quad (53)$$

Next we give the concentration result of $\tilde{R}_{P_n}(\mathbf{w}_{\text{pre}})$. Due to the definition of $\tilde{R}_{P_n}(\mathbf{w}_{\text{pre}})$, it can be rewritten as $R_{P_n^*}(\mathbf{w}_{\text{pre}})$ where P_n^* is the empirical distribution on $\{T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x}_i)\}$. Since $0 \leq f(\mathbf{w}_{\text{pre}}, \cdot) \leq M$ and $\{T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x}_i)\}$ are i.i.d draws from P_r^* . Azuma-Hoeffding's inequality (Corollary 2.20 in (Wainwright, 2019)) shows that with probability at least $1 - \theta$,

$$\tilde{R}_{P_n}(\mathbf{w}_{\text{pre}}) - R_{P_r^*}(\mathbf{w}_{\text{pre}}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{\text{pre}}, T_r^{\mathbf{w}_{\text{pre}}}(\mathbf{x}_i)) - R_{P_r^*}(\mathbf{w}_{\text{pre}}) \leq M \sqrt{\frac{\log(1/\theta)}{2n}}. \quad (54)$$

Hence we get our conclusion. \square

B.2. Proof of Theorem 5

With a little abuse of notation, let $r(2) = r/\epsilon_{\text{pre}}$ denoted by r in the proof, and $P_r^* \in \arg \max_{P \in B_{W_2}(P_0, r)} R_P(\mathbf{w})$. By Lemma 2, there exists $T_r^{\text{wpre}}(\mathbf{x}) \sim P_r^*$ with $\mathbf{x} \sim P_0$. Then we can find $\mathbf{z} \sim Q_0$ due to Kolmogorov’s Theorem. Let $T_r^{\text{wpre}}(\mathbf{z}) \sim Q_r^*$, we see

$$\begin{aligned} W_2(Q_0, Q_r^*)^2 &\leq \int_{\mathcal{X}} \|\mathbf{z} - T_r^{\text{wpre}}(\mathbf{z})\|^2 dQ_0(\mathbf{z}) \\ &\leq \int_{\mathcal{X}} \|\mathbf{z} - T_r^{\text{wpre}}(\mathbf{z})\|^2 |dQ_0(\mathbf{z}) - dP_0(\mathbf{z})| + \int_{\mathcal{X}} \|\mathbf{z} - T_r^{\text{wpre}}(\mathbf{z})\|^2 dP_0(\mathbf{z}) \\ &\leq D^2 \int_{\mathcal{X}} |dQ_0(\mathbf{z}) - dP_0(\mathbf{z})| + r^2 \\ &= 2D^2 \text{TV}(P_0, Q_0) + r^2. \end{aligned} \tag{55}$$

Thus $R_{Q_r^*}(\mathbf{w}_{\text{pre}}) \leq \epsilon_{\text{pre}}$. Similar to (52) and (53) we get the conclusion.

C. Hyperparameters

Table 4: Hyperparameters of adversarial training on CIFAR10.

Hyperparam	Std	Adv- ℓ_2	Adv- ℓ_∞
Learning Rate	0.1	0.1	0.1
Momentum	0.9	0.9	0.9
Batch Size	128	128	128
Weight Decay	5e-4	5e-4	5e-4
Epochs	200	200	200
Inner Loop Steps	-	8	8
Perturbation Size	-	2/12	2/255
Perturbation Step Size	-	1/24	1/510

Table 5: Hyperparameters of adversarial training on ImageNet.

Hyperparam	Std	Adv- ℓ_2	Adv- ℓ_∞
Learning Rate	0.1	0.1	0.1
Momentum	0.9	0.9	0.9
Batch Size	512	512	512
Weight Decay	5e-4	5e-4	5e-4
Epochs	100	100	100
Inner Loop Steps	-	3	3
Perturbation Size	-	0.25	2/255
Perturbation Step Size	-	0.05	1/510

Table 6: Hyperparameters of adversarial training on BERT base model.

Hyperparam	Std	Adv- ℓ_2	Adv- ℓ_∞
Learning Rate	3e-5	3e-5	3e-5
Batch Size	32	32	32
Weight Decay	0	0	0
Hidden Layer Dropout Rate	0.1	0.1	0.1
Attention Probability Dropout Rate	0.1	0.1	0.1
Max Epochs	10	10	10
Learning Rate Decay	Linear	Linear	Linear
Warmup Ratio	0	0	0
Inner Loop Steps	-	3	3
Perturbation Size	-	1.0	0.001
Perturbation Step Size	-	0.1	0.0005

D. Ablation Study

D.1. Effect of Perturbation Size

We study the effect of perturbation size r in adversarial training in bounds (5) and (6). We vary the perturbation size r in $\{2^{-5}/12, 2^{-4}/12, 2^{-3}/12, 2^{-2}/12, 2^{-1}/12, 2^0/12, 2^1/12, 2^2/12, 2^3/12, 2^4/12, 2^5/12, 2^6/12, 2^7/12\}$ for Adv- ℓ_2 and in $\{2^{-4}/255, 2^{-3}/255, 2^{-2}/255, 2^{-1}/255, 2^0/255, 2^1/255, 2^2/255, 2^3/255, 2^4/255\}$ for Adv- ℓ_∞ . The perturbation step size η_x in Algorithm 1 is set to be $r/4$ (Salman et al., 2020a). Experiments are conducted on CIFAR10 and the settings follow those in Section 5.1.1.

The results are shown in Figures 3 and 4. In the studied ranges, the accuracy on the OOD data from all categories exhibits similar trend, i.e., first increases and then decreases, as r increases. This is consistent with our discussion in Section 5.1.1 that there is an optimal perturbation size r for improving OOD generalization via adversarial training. For data corrupted

under types Fog, Bright and Contrast, adversarial training degenerates the performance in Table 1. We speculate this is because the three corruption types rescale the input pixel values to smaller values and the same perturbation size r leads to relatively large perturbation. Thus according to the discussion in Section 5.1.1 that there is an optimal r for improving OOD generalization, we suggest conducting adversarial training with a smaller perturbation size to defend these three types of corruption. Figures 3 and 4 also show that smaller optimal perturbation sizes have better performances for these three types of corruption.

D.2. Effect of the the Number of Training Samples

We study the effect of the number of training samples, as bounds (5) and (6) suggest that more training samples lead to better OOD generalization. We split CIFAR10 into 5 subsets, each of which has 10000, 20000, 30000, 40000 and 50000 training samples. The other settings follow those in Section 5.1.1. The results are in shown Figures 5 and 6.

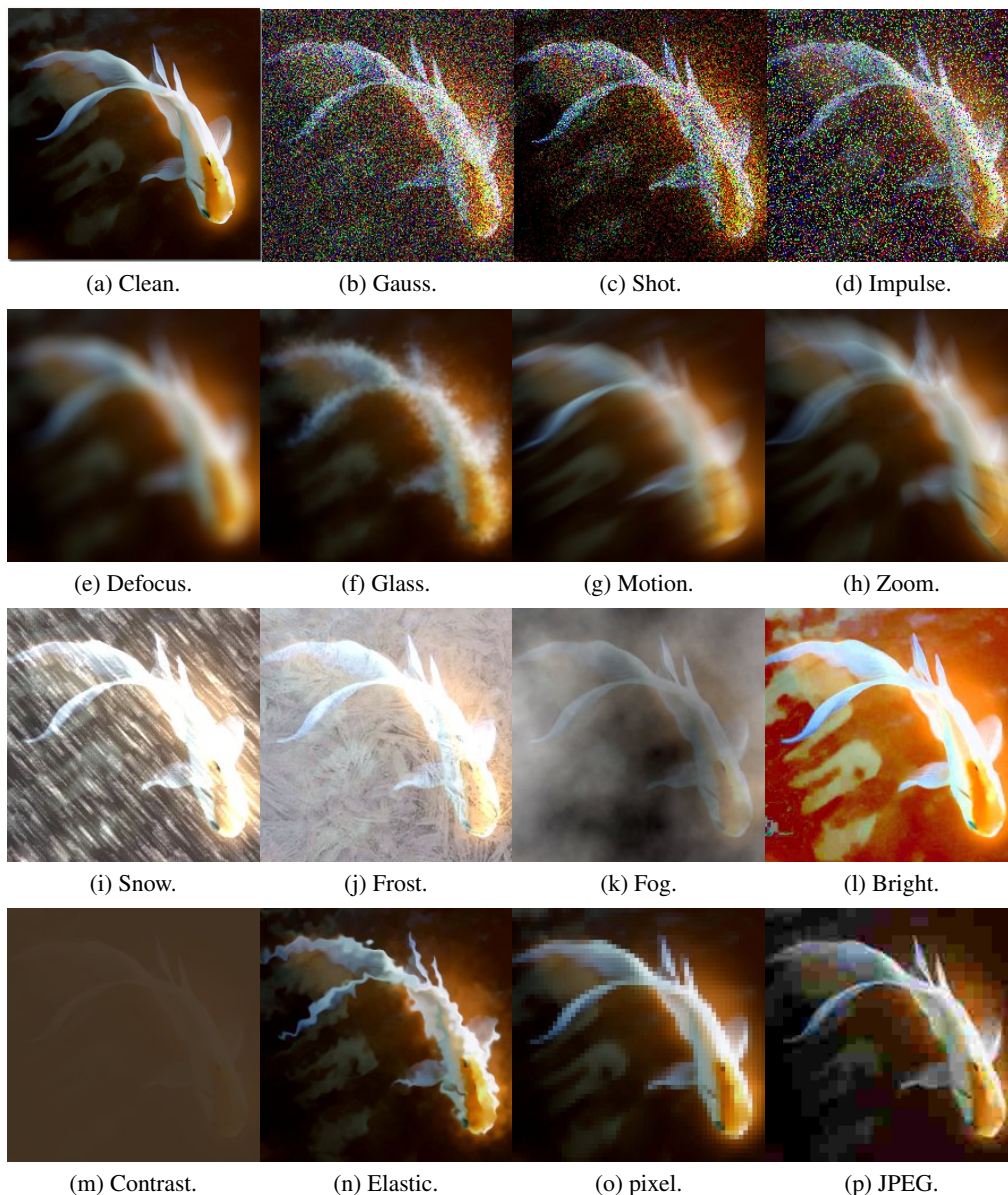


Figure 2: 15 types of artificially constructed corruptions from four categories: Noise, Blur, Weather, and Digital from the ImageNet-C dataset (Hendrycks & Dietterich, 2018). Each corruption has five levels of severity with figures under severity 5 are shown here.

Improved OOD Generalization via Adversarial Training and Pre-training

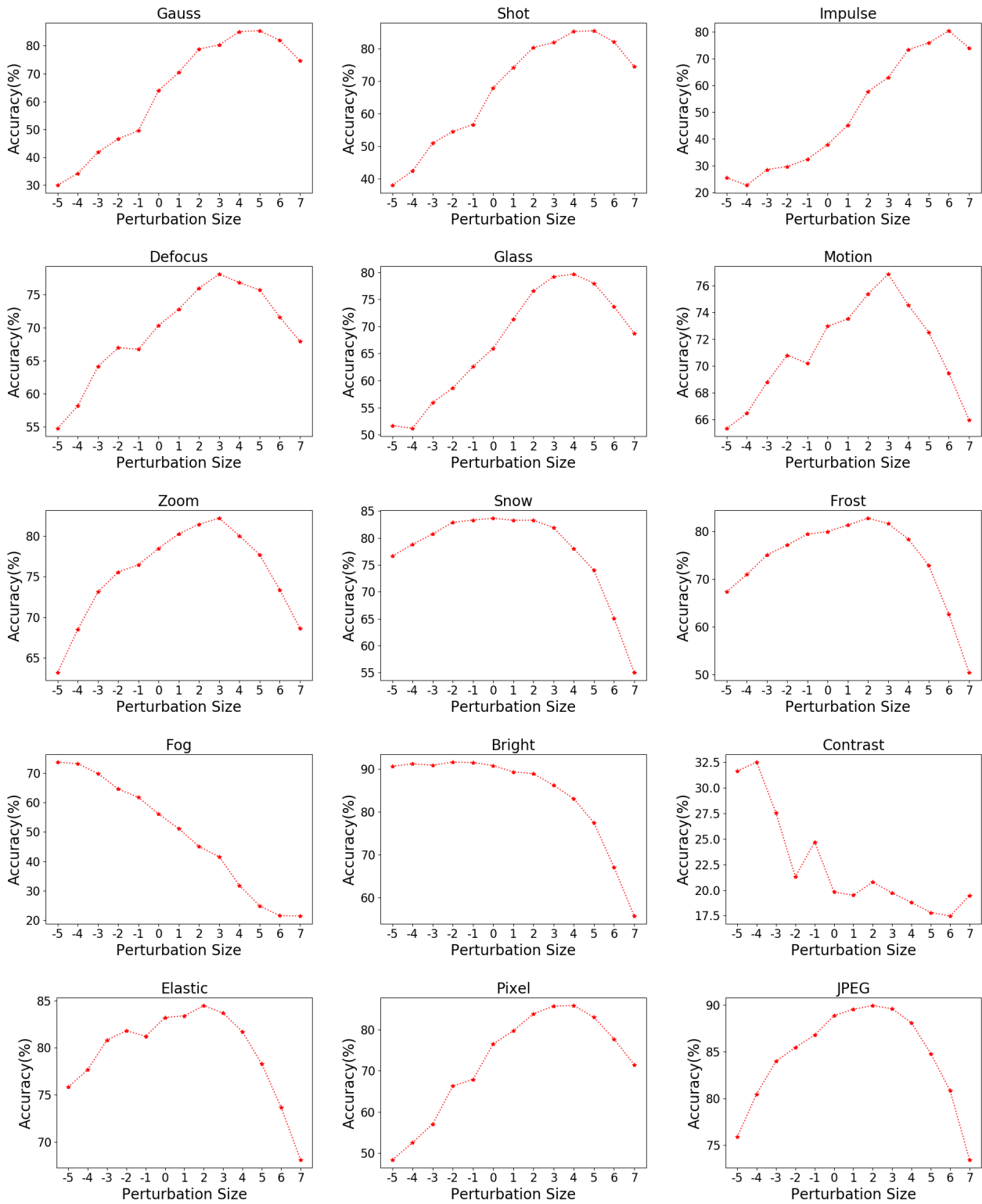


Figure 3: Accuracy of Adv- ℓ_2 on CIFAR10-C over various perturbation sizes. The x -axis means the perturbation size is $2^x/12$.

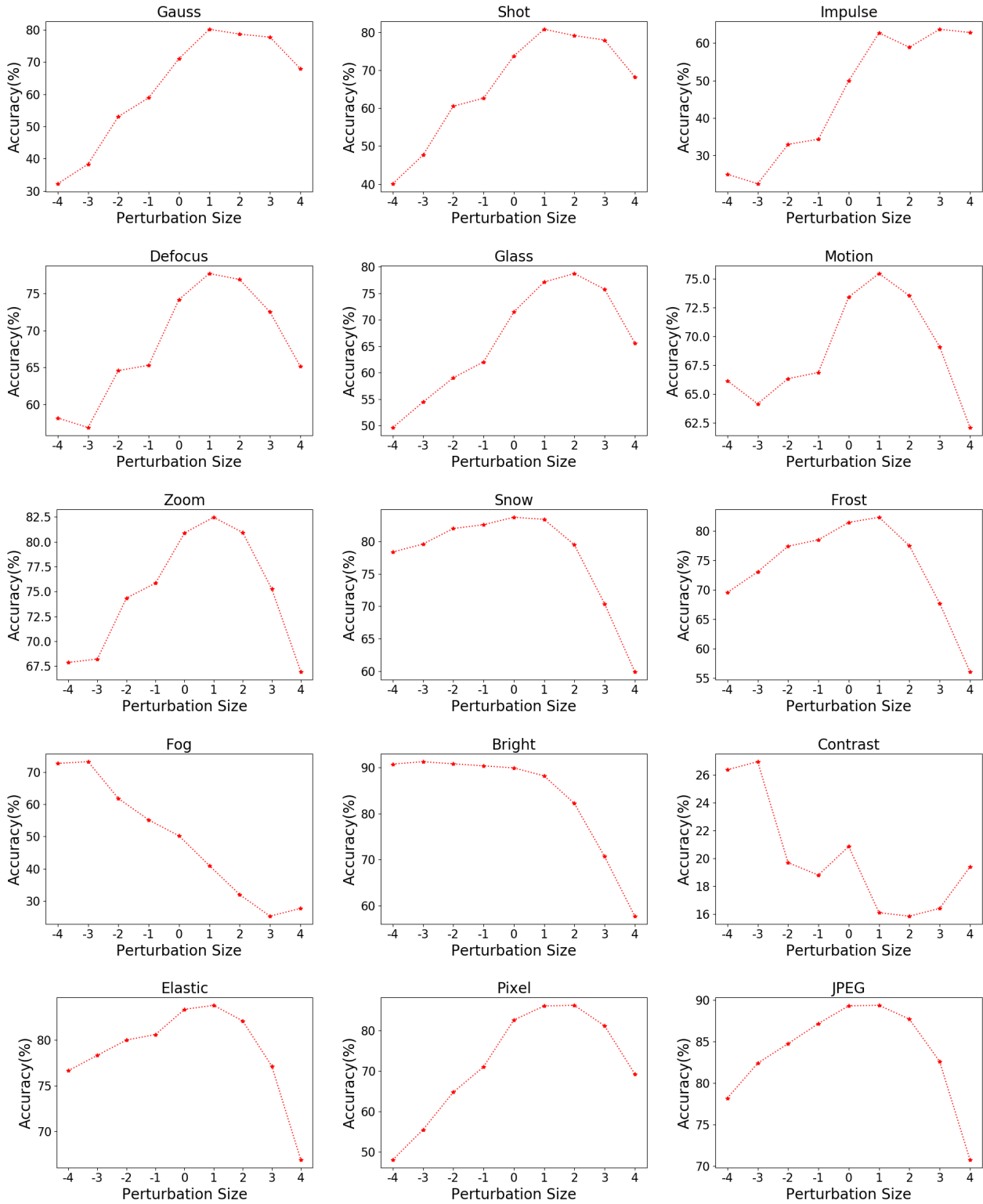


Figure 4: Accuracy of Adv- ℓ_∞ on CIFAR10-C over various perturbation sizes. The x -axis means the perturbation size is $2^x/255$.

Improved OOD Generalization via Adversarial Training and Pre-training

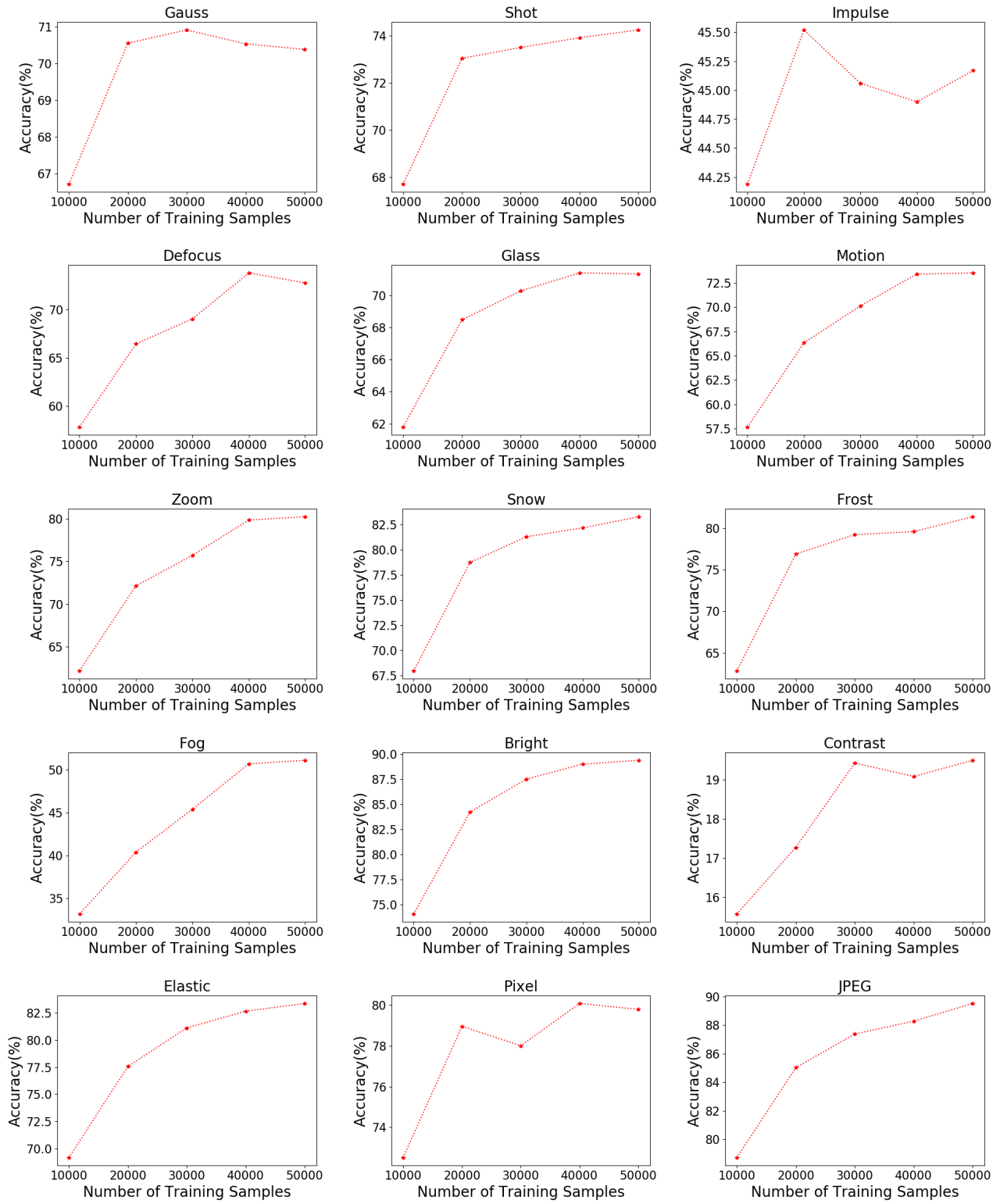


Figure 5: Accuracy of Adv- ℓ_2 on CIFAR10-C over various numbers of training samples.

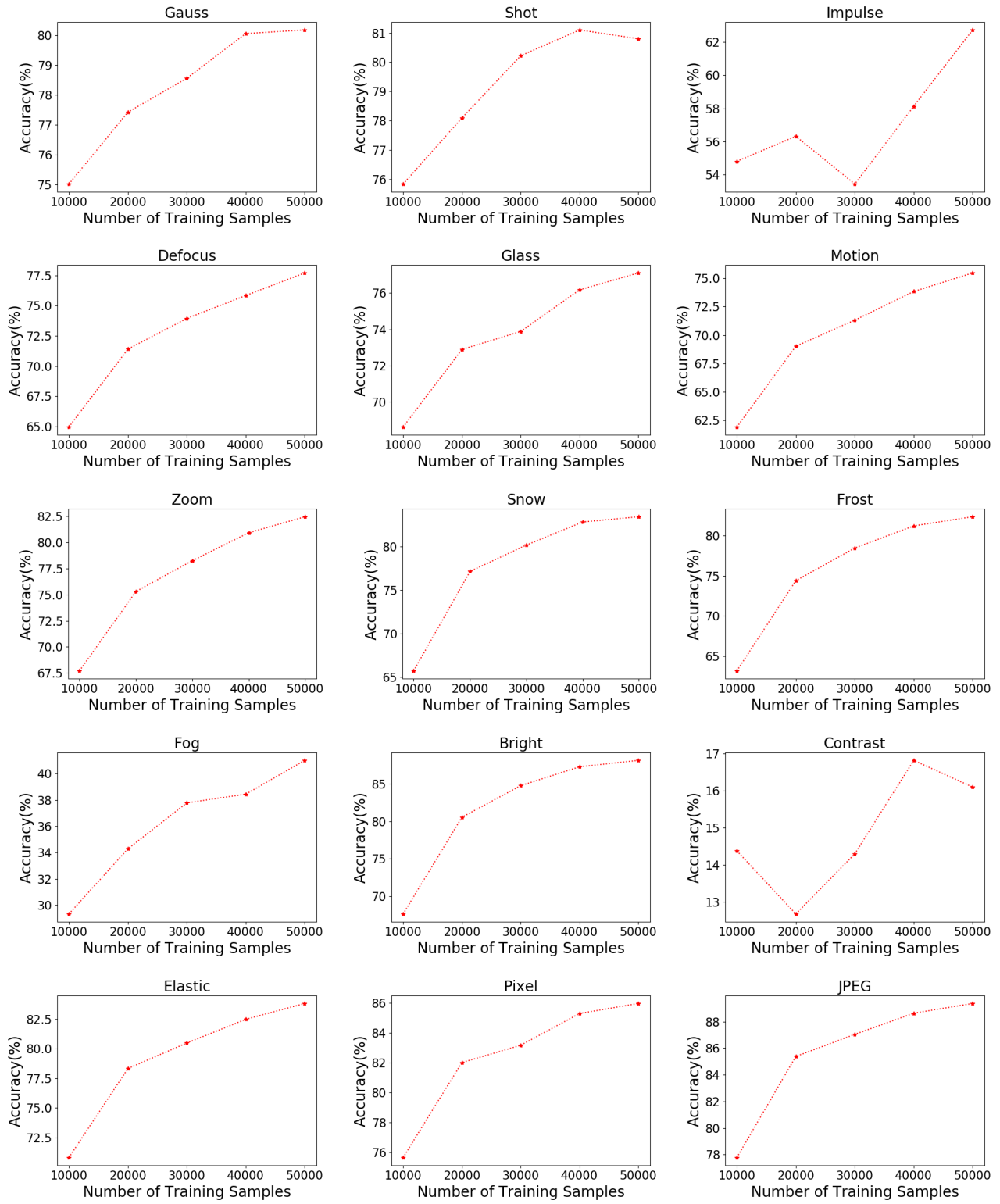


Figure 6: Accuracy of Adv- ℓ_∞ on CIFAR10-C over various numbers of training samples.