# 1 Test-Time Robustness

We now consider a different setting, focused on modeling corruptions at training time, rather than test time. For now assume that we are solving a classification task, so we wish to a predict a label $y \in \mathcal{Y}$ given covariates $x$ (for instance, $\mathcal{Y} = \{-1, +1\}$ corresponds to binary classification, or $\mathcal{Y} = \{0, \ldots, k-1\}$ to $k$-class classification). For a population distribution $p$, the setting is as follows:

- At train time, we observe samples $(x_1, \ldots, x_n) \sim p$.

- A test sample is generated by first drawing $(x, y) \sim p$, and replacing $x$ with $\bar{x}$ such that $d(x, \bar{x}) \le \epsilon$ for some distance $d$. We then observe $\bar{x}$ and hope to output $y$.

This can be placed into our previous framework by letting $\tilde{p} = p$ and saying that $D(p^*, \tilde{p}) \le \epsilon$ if there is a coupling $\pi$ from $\tilde{p}$ to $p^*$ such that $d(x, \bar{x}) \le \epsilon$ and $y = y'$ almost surely for $(x, y, \bar{x}, y') \sim \pi$. In the language of Wasserstein distance, this says that $W_{d,\infty}(\tilde{p}, p^*) \le \epsilon$ (here we define $W_{c,k}(p, q) = \inf_{\pi \in \Pi(p,q)} \mathbb{E}_{(z,z') \sim \pi}[c(z, z')^k]^{1/k}$).

However, there are some important differences from the setting from before:

- We generally think of $p = \tilde{p}$ as the "nice" distribution and $p^*$ as the "ugly" distribution.

- Since $p^*$ is the "ugly" distribution, we make no distributional assumptions $\mathcal{G}$ about $p^*$.

- On the other hand, we also generally consider much "smaller" perturbations than in the previous case. For instance, $d(x, \bar{x}) \le \epsilon$ should imply that $x$ and $\bar{x}$ are close enough that they have the same label $y$.

A final point is that the worst-case test loss is directly observable, at least given infinitely many samples from $p$. Indeed, if the (non-robust) loss is

$$L_0(p, \theta) = \mathbb{E}_{(x,y) \sim p}[\ell(\theta; x, y)], \tag{1}$$

then the robust loss is

$$L(p, \theta) = \mathbb{E}_{(x,y) \sim p}\left[ \sup_{\bar{x}:d(x,\bar{x}) \le \epsilon} \ell(\theta; \bar{x}, y) \right]. \tag{2}$$

(The worst-case test loss was also observable for train-time corruptions, but it involved a much more complex supremum and depended on $\mathcal{G}$.)

Consequently, given samples $(x_1, y_1), \ldots, (x_n, y_n)$, a natural estimator is

$$\hat{\theta} = \arg\min_\theta \rho(\theta) + \frac{1}{n} \sum_{i=1}^n \sup_{d(x_i, \bar{x}_i) \le \epsilon} \ell(\theta; \bar{x}_i, y_i), \tag{3}$$
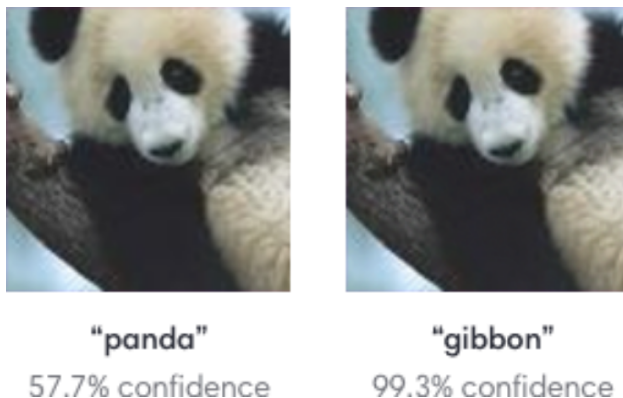
where $\rho(\theta)$ is a regularizer and the other term is an empirical estimate of the worst-case loss. Indeed, this is essentially the estimator we will consider throughout. However, there are several issues to address:

- How well does $\hat{\theta}$ generalize from the training points to the population distribution?

- The supremum over $d(x, \bar{x}) \le \epsilon$ is generally computationally intractable. What happens if we only approximate this supremum, and what are good approximation schemes?

- How robust are results to the choice of $d$ itself?

Answering these questions will involve a mix of empirical data and theoretical analysis.

## 1.1 Getting Oriented: Examples of $d$ and basic empirical results

A basic motivating example takes $x$ to be an image containing some object $y$. Thus for instance $x \in [0,1]^{224 \times 224 \times 3}$ is a 150528-dimensional vector (corresponding to a $224 \times 224$ RGB image), and $y$ is one of $k$ different object classes (in the most popular computer vision benchmark, ImageNet, we have $k = 1000$). We then take $d(x, \bar{x}) = \|x - \bar{x}\|_\infty$, meaning that we are allowed to perturb each pixel by at most $\epsilon$. A striking phenomenon is that regularly-trained machine learning models are often non-robust at even very small values of $\epsilon$. For instance, the following two images have different classifications under a neural network:
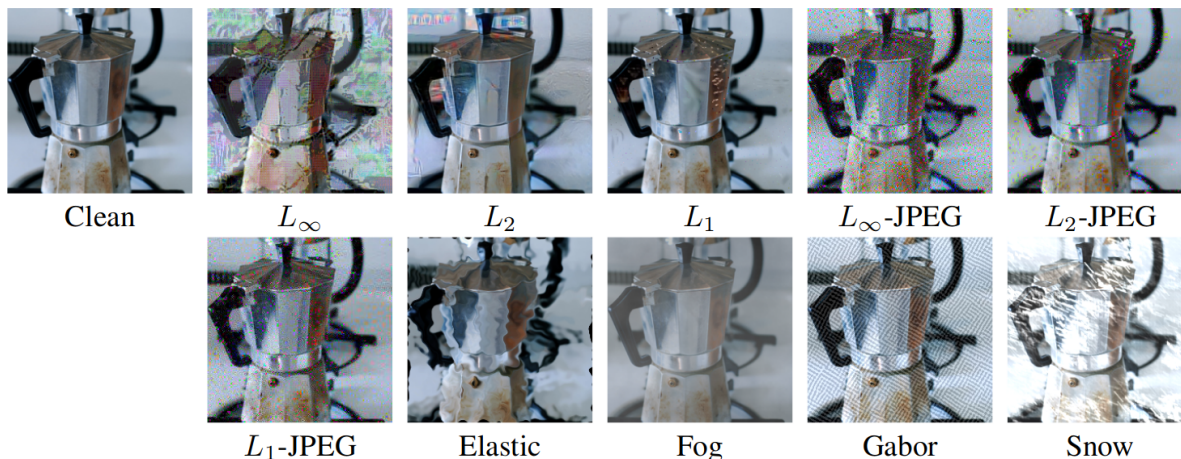


"panda"
57.7% confidence

"gibbon"
99.3% confidence

This corresponds to a perturbation in $\ell_\infty$ of $\epsilon = 0.007$. This is not an isolated phenomenon, and many machine learning models can be successfully adversarially perturbed at small $\epsilon$ for most of their inputs.

We could consider many perturbation measures beyond $\ell_\infty$:

- $d(x, \bar{x}) = \|x\bar{x}\|_\infty$: the maximum amount any pixel changes

- $d(x, \bar{x}) = \|x - \bar{x}\|_1$: the total amount that the pixels change

- $d(x, \bar{x}) = \|\mathsf{JPEG}(x) - \mathsf{JPEG}(\bar{x})\|_2$: the $\ell_2$-distance in wavelet space

- The amount of "stretch" needed to move $x$ to $\bar{x}$

- The amount of "fog", "snow", etc. needed to move $x$ to $\bar{x}$.

These and others are illustrated below:



| Clean | $L_\infty$ | $L_2$ | $L_1$ | $L_\infty$-JPEG | $L_2$-JPEG |

| $L_1$-JPEG | Elastic | Fog | Gabor | Snow |

Note that for many of these, $d$ is not a metric (it isn't even symmetric, e.g. for snow).

There are four important phenomena that we will elaborate on below:

- Adversarial perturbations have a much larger effect than random perturbations.

- Adversarial robustness is difficult to measure and most papers get it wrong enough that the results are meaningless.

- Robustness for one choice of $d$ only partially generalizes to other choices of $d$.

- Minimizing the robust loss seems to increase overfitting for neural networks.

**Worst-case vs. random perturbations.** As a simplified case suppose that we are doing binary classification and our model is linear: $f_\theta(x) = \text{sign}(\langle \theta, x \rangle)$. We then consider two types of perturbations:

- Random Gaussian perturbations: $\bar{x} = x + z$, where $z \sim \mathcal{N}(0, \sigma^2 I)$ (standard deviation $\sigma$ per coordinate).

- Worst-case $\ell_2$-perturbations: $\bar{x} = x + z$, where $\|z\|_2 \leq \sigma\sqrt{d}$.

Note that the $\ell_2$-norm of the perturbation is approximately the same in both cases. For a point $x$, define the *margin* $\gamma_x = |\langle \theta, x \rangle|/\|\theta\|_2$. How big does the margin need to be for the output $f_\theta$ to be robust to perturbations?

For the Gaussian perturbations, we have $\langle \theta, \bar{x} \rangle = \langle \theta, x \rangle + \sigma\|\theta\|_2 \mathcal{N}(0, 1)$, so the output is robust as long as $\gamma_x \gg \sigma$. However, for worst-case $\ell_2$-perturbations we have $\langle \theta, \bar{x} \rangle = \langle \theta, x \rangle + \langle \theta, z \rangle$, which in the worst-case is $\langle \theta, x \rangle \pm \sigma\sqrt{d}\|\theta\|_2$. Thus we are only robust to worst-case perturbations when $\gamma_x > \sigma\sqrt{d}$. There is a $\sqrt{d}$ factor difference between random and worst-case perturbations! In fact, a Gaussian perturbation with *per-coordinate* magnitude $\sigma$ is most similar to a worst-case perturbation with *overall* magnitude $\sigma$.

**Adversarial robustness is difficult to measure.** To measure adversarial robustness, we seek to compute

$$\ell_{\text{robust}}(\theta; x, y) = \sup_{\bar{x}: d(x, \bar{x}) \leq \epsilon} \ell(\theta; \bar{x}, y) \tag{4}$$

for a given $x$ and $y$, which is computationally intractable in general. For instance, $\ell$ is often the loss function for a neural network, which is highly non-convex not just in $\theta$ but also in $x$ (which is the relevant variable for our purposes). One idea would be to approximate the supremum via some local search algorithm (e.g. projected gradient descent). This gives us some proxy loss $\ell_{\text{proxy}} \leq \ell_{\text{robust}}$.

The issue arises when we take the natural step of optimizing $\ell_{\text{proxy}}$ to try to find a robust model $\theta$. Since $\ell_{\text{proxy}}$ underestimates $\ell_{\text{robust}}$, this highly incentivizes finding $\theta$ where $\ell_{\text{proxy}}$ is a bad approximation to $\ell_{\text{robust}}$! For instance, if our proxy loss is obtained via gradient descent, it incentivizes $\theta$ for which gradient descent on $\bar{x}$ doesn't converge well. Even if we don't explicitly optimize $\ell_{\text{proxy}}$, if we design models with $\ell_{\text{proxy}}$ in mind we as the researcher might accidentally find ways to make $\ell_{\text{proxy}}$ but not $\ell_{\text{robust}}$ small. And since $\ell_{\text{proxy}}$ is the only thing we measure, we have no straightforward way of even noticing this!

I estimate that around 95% of papers on test-time robustness succumb to this issue, where the proxy loss measured in the paper is small but the true robust loss is large. Some good papers that discuss this are **?** and **?**, which both also provide best practices for evaluating robustness. (See also **?** for an earlier paper making similar points for the detection of adversarial perturbations.) Some example best practices are:
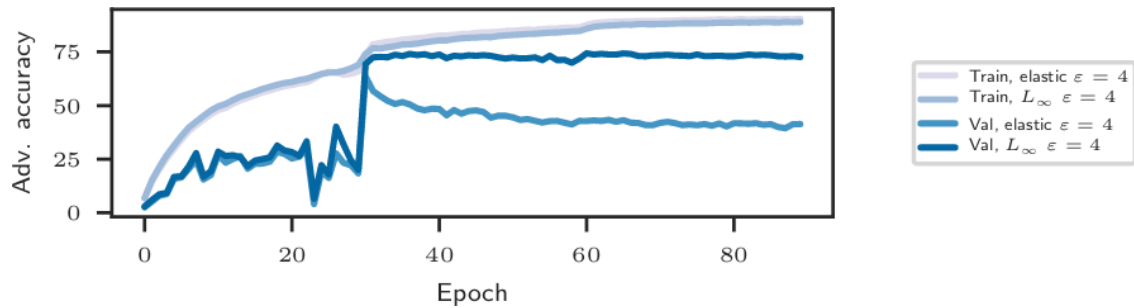
- If your proxy is based on some optimization algorithm, give the optimizer much more power at evaluation time than at training time (e.g. use 10 gradient descent steps at train time but 1000 at test time).

- Ensure that your proxy works well by making sure that when $\epsilon$ is large enough that accuracy should be zero, the proxy also gives an accuracy of zero.

- Make sure that other correlated measures (such as robustness to random perturbations) are also improving. This is more controversial because not everyone agrees on what measures should correlate with adversarial robustness.

**Robustness to choice of $d$.** While this is a nuanced story, the zeroth order summary is that robust for fairly different choices of $d$ is positively correlated at low $\epsilon$, but eventually becomes negatively correlated at high $\epsilon$ (although aiming for robustness under the wrong $d$ often still gives better results than doing nothing). Here is a grid illustrating this:

Here each row is a model that was trained to be robust to a certain type of perturbation, and each column is a measure of the robustness against a perturbation (higher is better). We will discuss in the next section how to actually perform this training. We see that even for fairly perturbations as different as elastic warping and snow, training against one perturbation helps with the other one. See **?** for a more detailed investigation of this.

**Robustness hurts generalization.** Training a model to be robust, especially at large $\epsilon$, can harm generalization performance. This becomes even more pronounced if we allow multiple perturbations, for instance if $d$ allows both elastic warping and $\ell_\infty$ perturbations, we see overfitting even at a moderately-sized $\epsilon = 4$ (this is $4/255$ so approximately 0.016):



We see similar phenomena for a single perturbation, but need $\epsilon$ to be larger (e.g. around $32/255$ for $\ell_\infty$, which is large enough to be hard for humans).