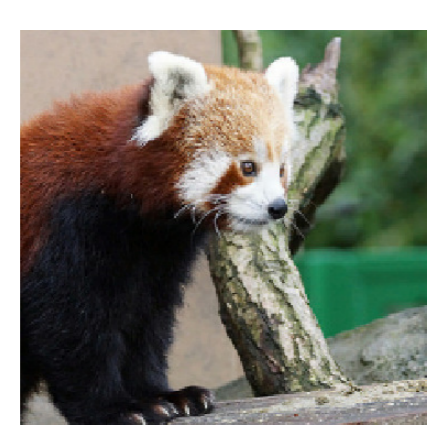


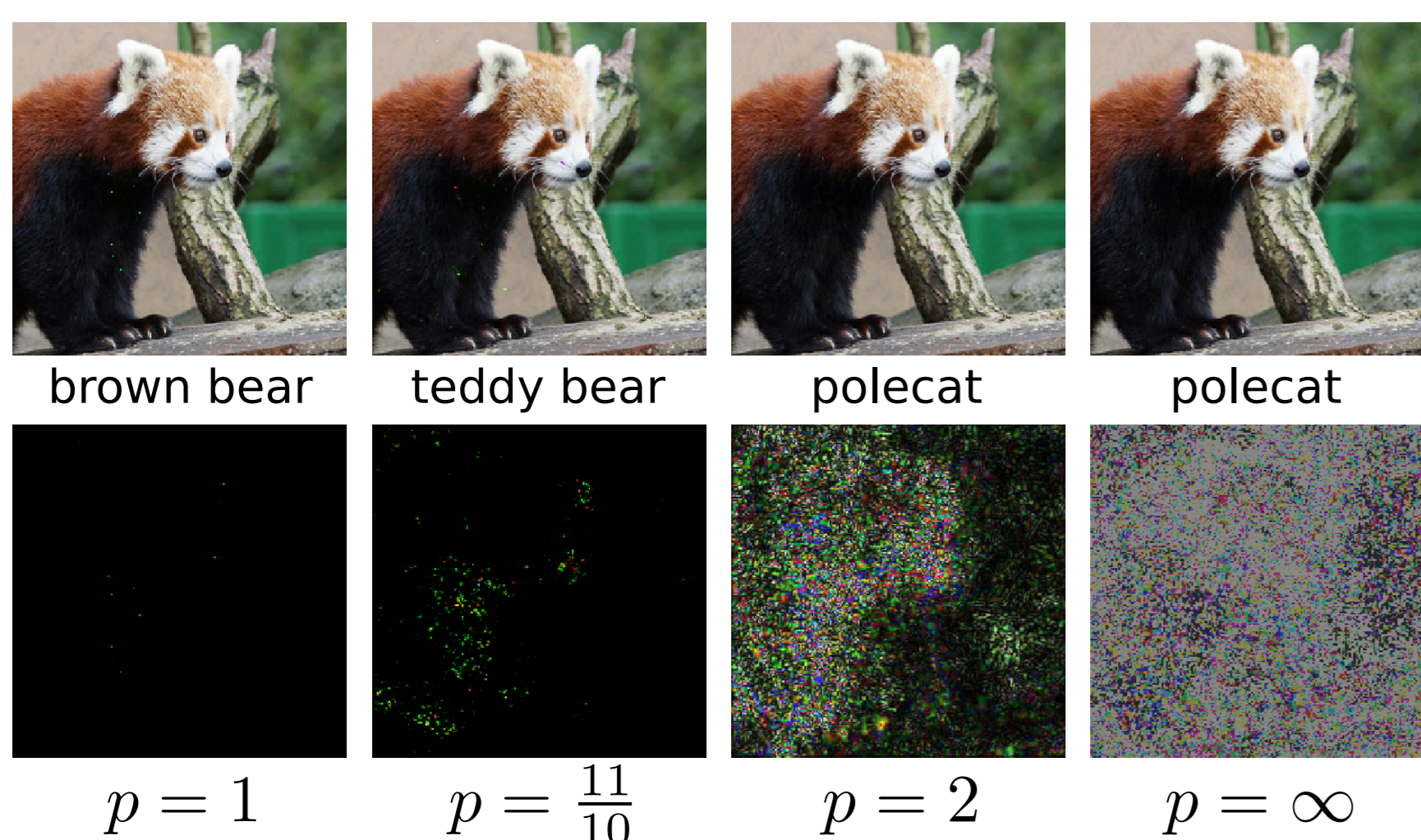
### 1. Adversarial and Random Perturbations

- Goal: **Relate robustnesses to adversarial and random noise.**
- Classifier  $f: \mathbb{R}^d \rightarrow \mathbb{R}^L$ ,  $g = \operatorname{argmax} f$ .
- Adversarial perturbation w.r.t. norm  $\ell_p$ :

$$r_p^*(x) = \operatorname{argmin}_r \{ \|r\|_p \text{ s.t. } g(x+r) \neq g(x) \}.$$



red panda  
(unperturbed)



**Figure 1:** Adversarial perturbations of an image for VGG-19 with different  $p$ -norms.

- Robustness to random perturbation, for  $v \sim \nu$ :  
$$r_{\nu,\varepsilon}(x) = \min_{\alpha} \{ |\alpha| \text{ s.t. } \mathbb{P}_v \{ g(x + \alpha v) \neq g(x) \} \geq \varepsilon \}.$$
- Goal: **Derive lower and upper bounds as well as an estimate on**  $\frac{r_{\nu,\varepsilon}(x)}{\|r_p^*(x)\|_p}$ .

### 2. Bounds on Linear Classifiers

$$f(x) = w^T x + b$$

#### 2.1 Uniformly Distributed Noise in the $\ell_p$ Ball

**Theorem 1.** Let  $p, p' \in [1, \infty]$  such that  $\frac{1}{p} + \frac{1}{p'} = 1$ . Then, for  $\varepsilon$  small enough:

$$\zeta_1(\varepsilon) d^{1/p} \frac{\|w\|_{p'}}{\|w\|_2} \leq \frac{r_{p,\varepsilon}(x)}{\|r_p^*(x)\|_p} \leq \zeta_2(\varepsilon) d^{1/p} \frac{\|w\|_{p'}}{\|w\|_2}.$$

- Depends on the choice of orthonormal basis if  $p \neq 2$ .
- For a typical  $w$ , each bound is of the form:

$$C(\varepsilon, p) \sqrt{d}.$$

**Proof sketch of simple special case. Lower bound for  $p = \infty$ :**

- $r_p^*(x) = \frac{|f(x)|}{\|w\|_1}$  and  $\mathbb{P}_v \{ g(x + \alpha v) \neq g(x) \} = \mathbb{P}_v \left\{ |\alpha| \sum_{i=1}^d w_i v_i \geq \|w\|_1 |f(x)| \right\}.$
- $(v_i)_i$  are i.i.d. and uniform over  $[-1, 1]$ : can apply Hoeffding and derive a lower bound.

**Ideas for the general case.**

- Lower bound: use Markov's inequality and estimates of  $\mathbb{E}_v \left[ (w^T v)^k \right]$ .
- Upper bound: use Paley-Zygmund's inequality and the previous estimates.

### 2.2 Gaussian Noise

**Theorem 2.** Let  $\Sigma$  be a  $d \times d$  positive semidefinite matrix. Then, for  $\varepsilon$  small enough:

$$\zeta'_1(\varepsilon) \frac{\|w\|_2}{\|\sqrt{\Sigma} w\|_2} \leq \frac{r_{\Sigma,\varepsilon}(x)}{\|r_2^*(x)\|_2} \leq \zeta'_2(\varepsilon) \frac{\|w\|_2}{\|\sqrt{\Sigma} w\|_2}.$$

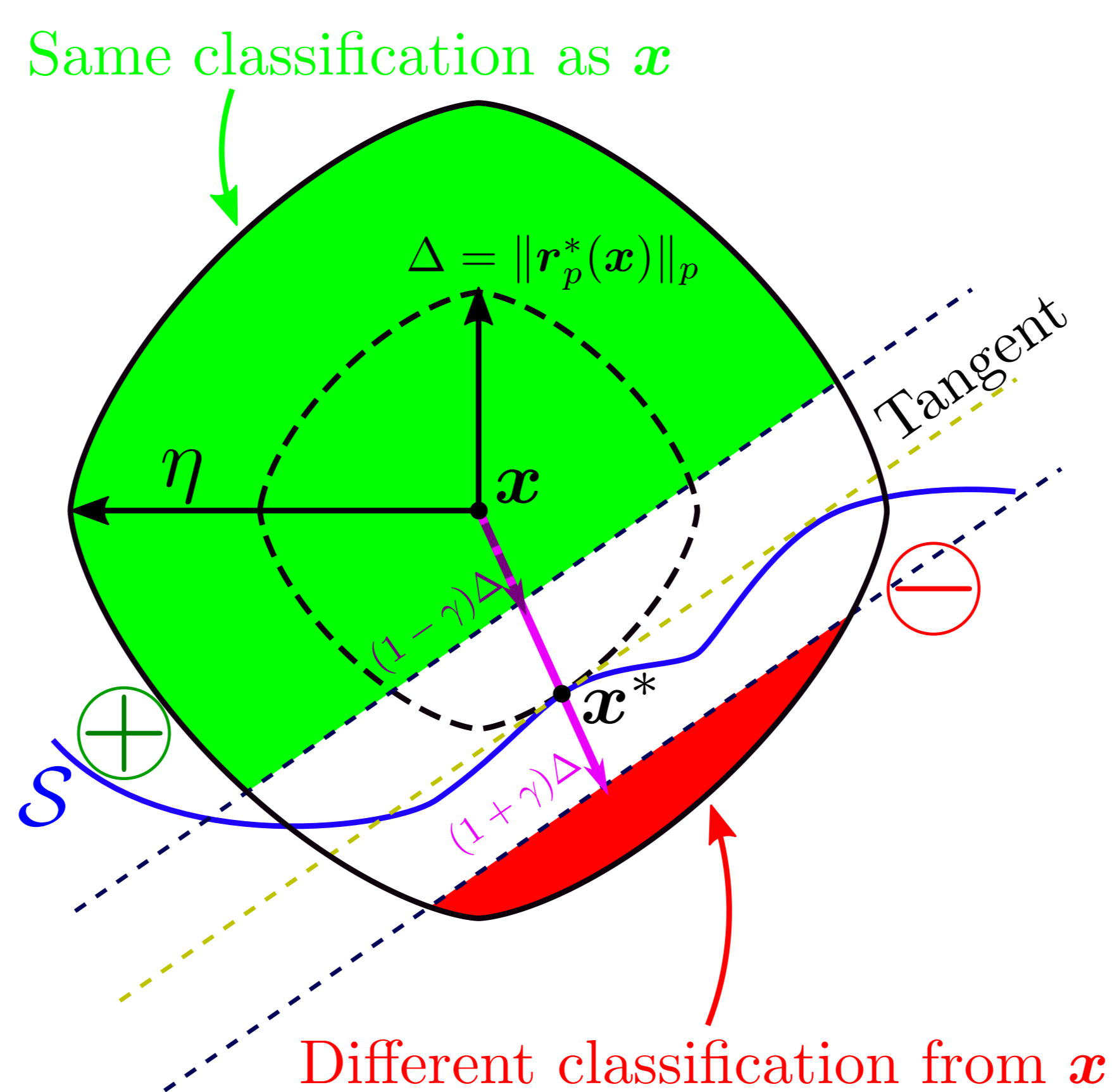
- For a typical  $w$ , each bound is of the form:

$$C'(\varepsilon, p) \sqrt{d}.$$

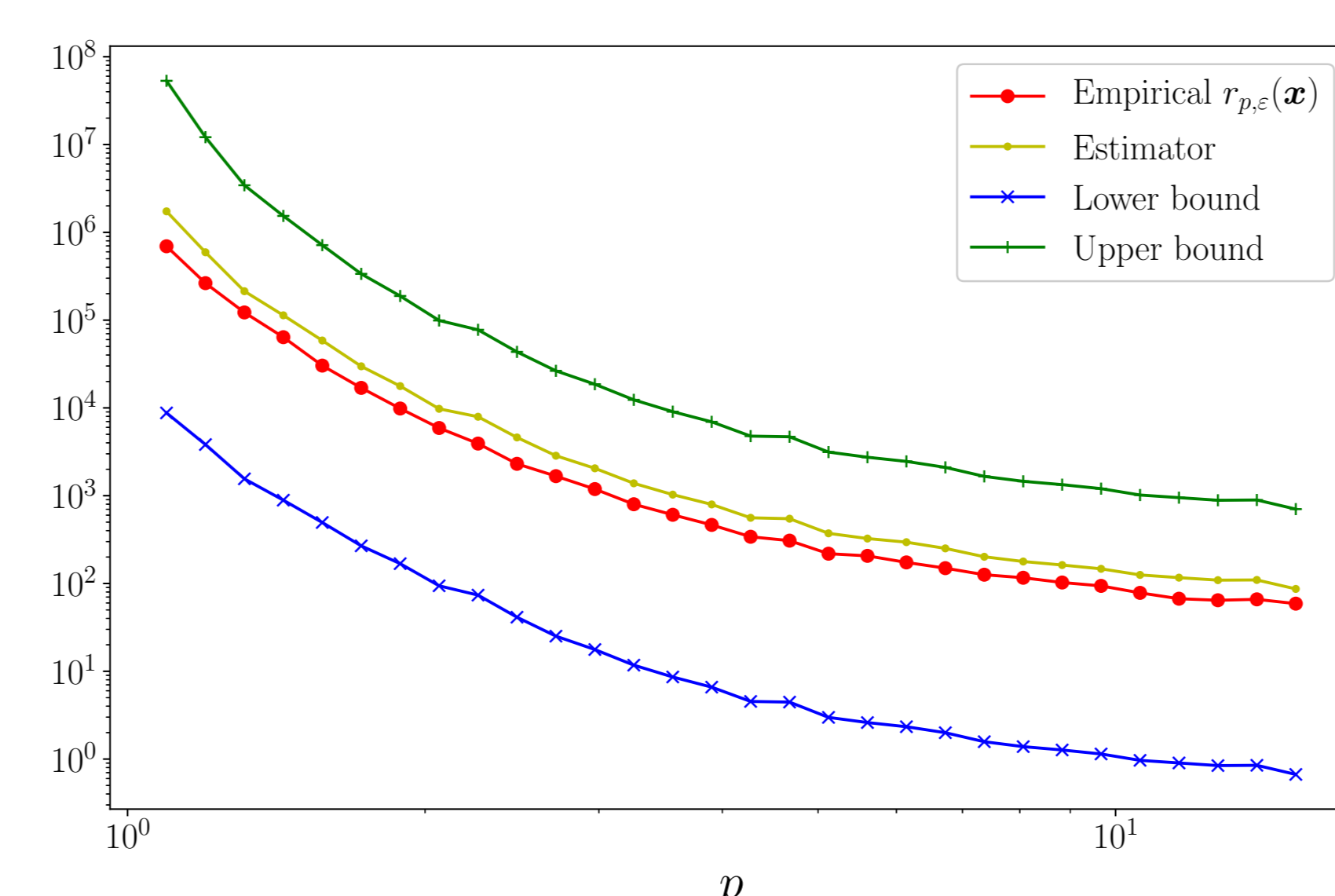
- $\Sigma$  may depend on  $x$ .

### 3. Extension to Locally Approximately Flat (LAF) Classifiers

- A classifier is LAF at some point  $x$  if the decision boundary can be approached by a plane in a given ball centered at  $x$ .
- In this setting, the results on linear classifiers are still valid, up to some constants in the bounds.
- The normal vector  $w$  may be naturally replaced by the gradient of  $f$  at the closest point  $x^*$  on the decision boundary.
- Experiments indicate that the LAF assumption is reasonable for state-of-the-art deep neural networks.



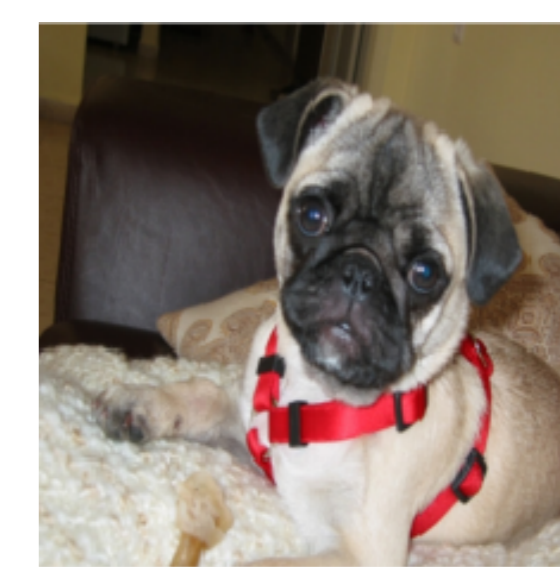
**Figure 2:** Illustration of the LAF model.



**Figure 3:** Experiments on a deep neural network (VGG-19, ImageNet dataset).

### 4. Applications: Robustness to Quantization

- Image quantization (discretization of the range of pixel values):  $x \mapsto Q(x)$ .
- Can assume  $Q(x) \sim \mathcal{U}(\mathcal{B}_\infty(x, \frac{\Delta}{2}))$ .
- $\Delta$ : quantization step size.
- $L_q = \frac{255}{\Delta}$ : number of quantization levels corresponding to  $\Delta$ .
- Our results allow us to estimate the level of quantization needed so that a quantized image is still classified correctly.



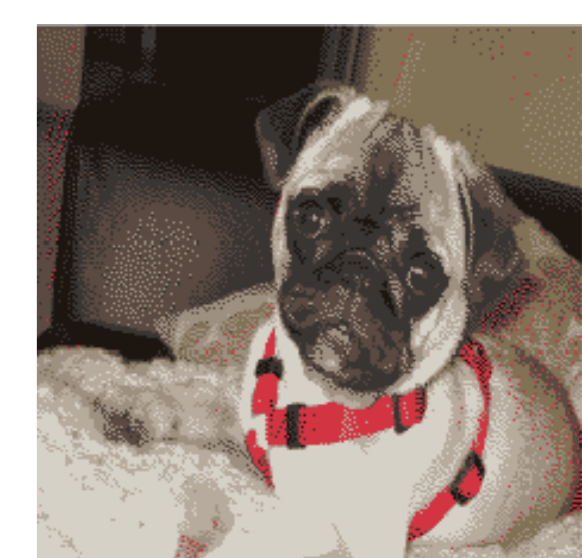
pug  
Original



mail  
1 bit

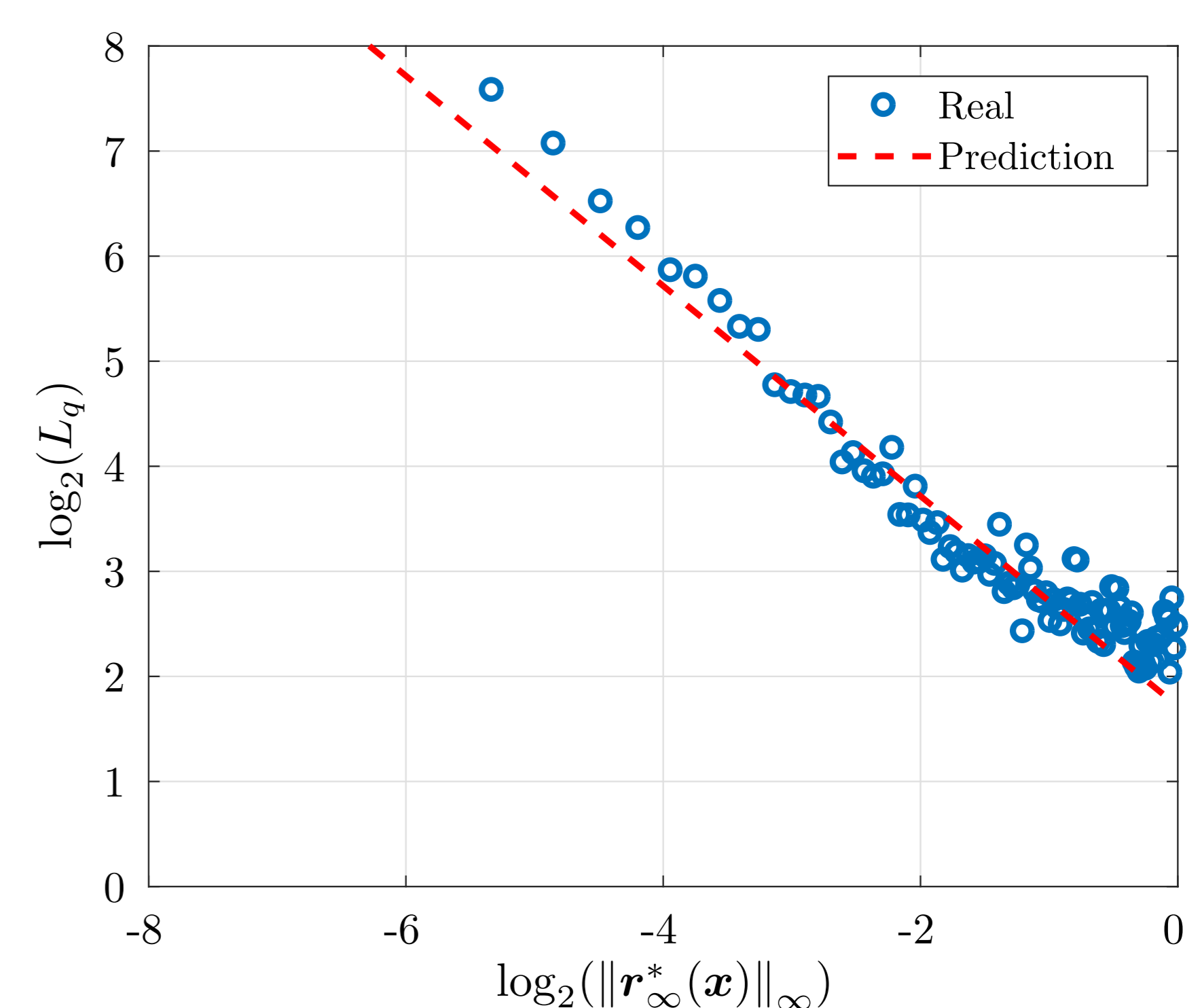


bulldog  
2 bits



pug  
3 bits

**Figure 4:** Illustration of the effects of quantization.



**Figure 5:** Experiments on quantization (VGG-19, ImageNet dataset).

### References

- Barthe, F., Guédon, O., Mendelson, S., and Naor, A. (2005). A probabilistic approach to the geometry of the  $\ell_p^n$ -ball. *The Annals of Probability*, 33(2).
- Fawzi, A., Moosavi-Dezfooli, S., and Frossard, P. (2016). Robustness of classifiers: from adversarial to random noise. In *Advances in Neural Information Processing Systems (NIPS)* 29.
- Moosavi-Dezfooli, S.-M., Fawzi, A., and Frossard, P. (2016). Deepfool: a simple and accurate method to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.