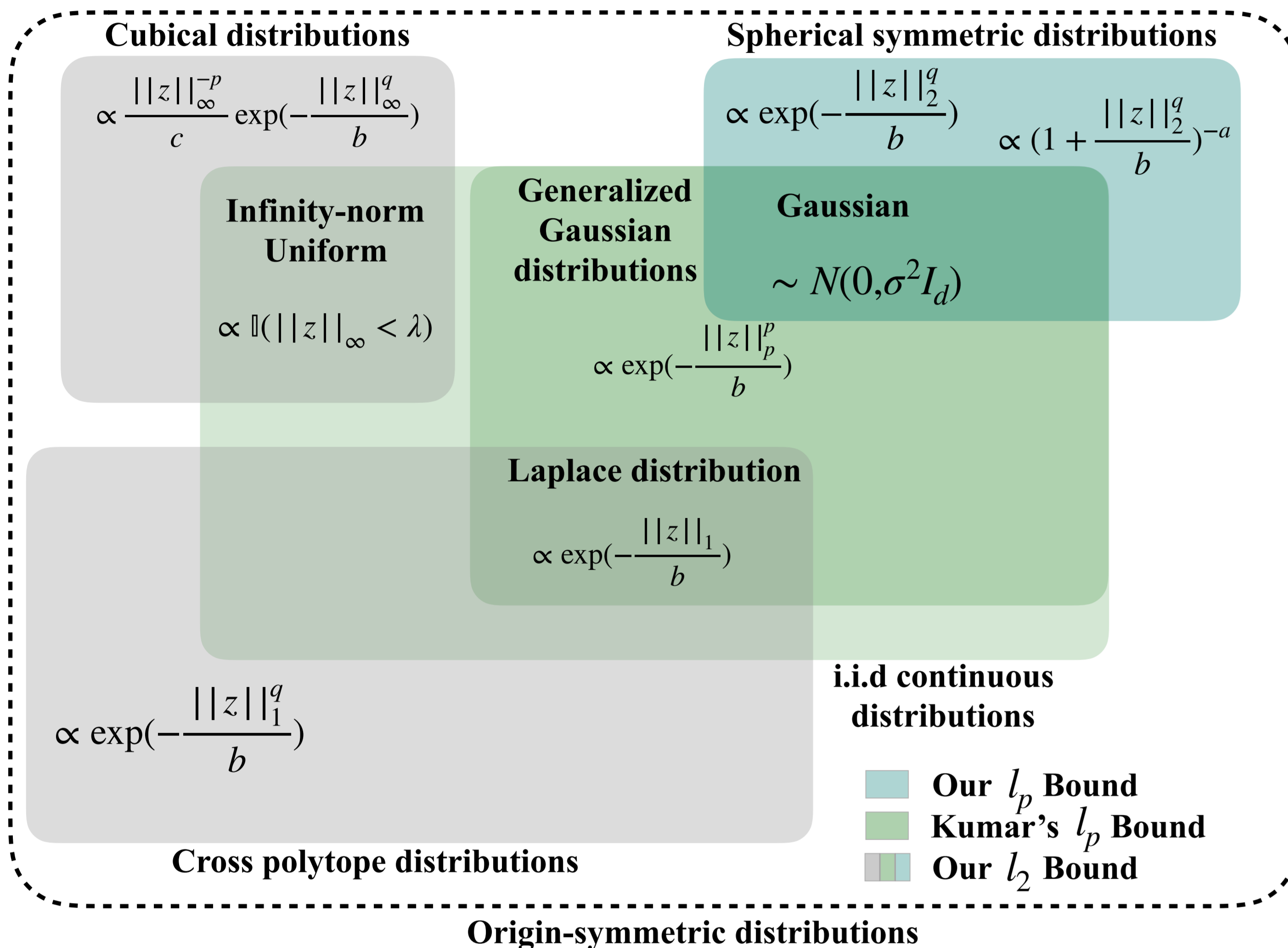


TL;DR: We derived two upper bounds on the certified radius with randomized smoothing method for origin-symmetric distributions. The radius scales as  $O(1/\sqrt{d})$ .



## Main Research Question

What is the best result we can hope for with randomized smoothing?

## Randomized smoothing is important

Randomized smoothing is currently the only SOTA certificate that scales to large networks and different settings.

**Smoothed classifier:** given a base classifier  $f$  and a sample  $x$  with label  $c$ , the smoothed classifier  $g$  with distribution  $q$  is

$$g_f(x) := E_{z \sim q}[1_{f(x+z)=c}]$$

**Certified radius:** the largest  $r$  such that for any perturbation  $\delta$  with norm  $\|\delta\|_2 < r$ , the prediction of perturbed samples does not change.

**Certified radius in randomized smoothing:** We calculate a **lower bound** of certified radius of the smoothed classifier by finding a lower bound of  $g_f(x + \delta)$ .

## Methods to find the upper bound

**Functional optimization method:**

Provide a tight lower bound of  $g_f(x + \delta)$ .

Optimization problem:

$$\begin{aligned} \underline{g}(x + \delta) &:= \min_h g_h(x + \delta) \\ \text{s.t. } g_h(x) &= g_f(x) \end{aligned}$$

where  $h$  is in the feasible set of all classifiers.

The certified radius is the largest norm of  $\delta$  such that

$$\underline{g}(x + \delta) > 0.5$$

**Intuition:** we can find an upper bound of certified radius by selecting a classifier  $h$  and perturbation  $\delta$  such that  $g_h(x) = g_f(x)$  and  $g_h(x + \delta) < 0.5$ , then

$$\underline{g}(x + \delta) < g_h(x + \delta) < 0.5$$

and the certified radius is upper bounded by the norm of  $\delta$ .

We need to find a classifier and a perturbation that simplify the calculation.

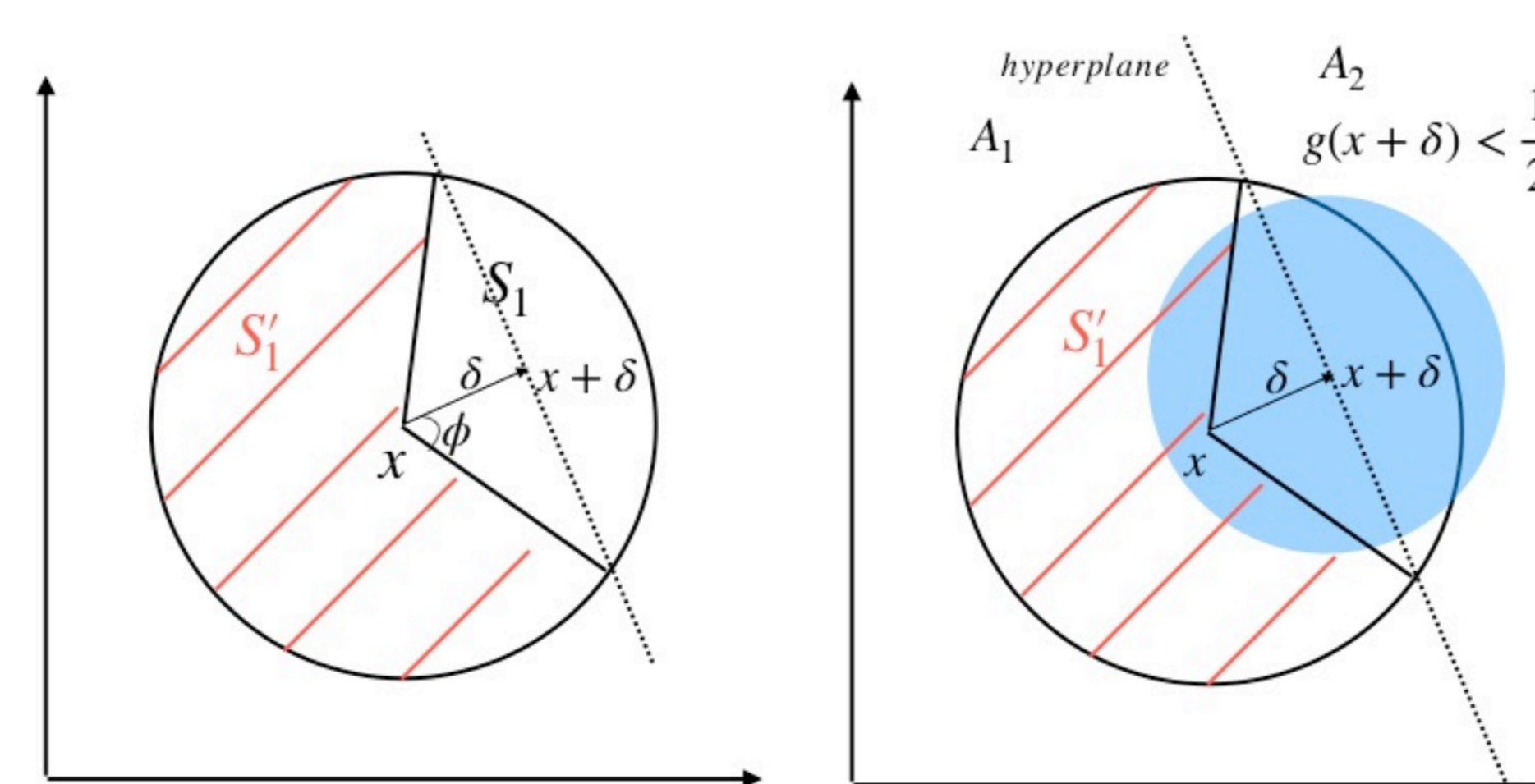
## Bound the $\ell_2$ radius

$S_1$ : a hyperspherical sector of a  $d$ -dimensional ball.

The classifier we defined is  $1_{\{x \in S_1\}}$  and  $\delta$  is orthogonal to the hyperplane.

**1<sup>st</sup> upper bound:**

$$r < \|\delta\|_2 < \frac{5}{\sqrt{d}} \Psi^{-1}\left(\frac{g(x)}{1 - 5 \times 10^{-7}}; q\right), \quad \Psi(R; q) = \int_{B_R} q(z) dz$$

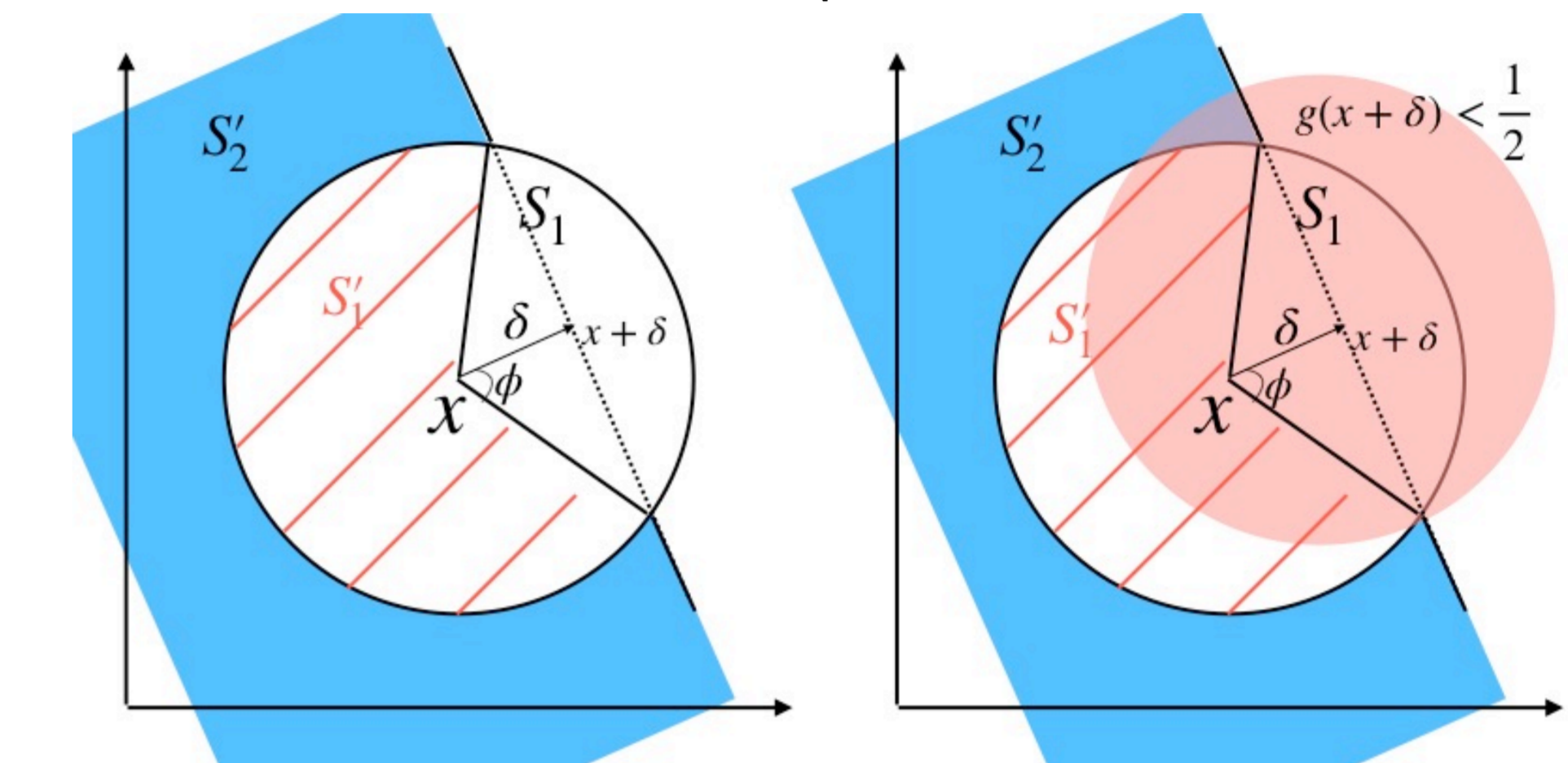


## Alternative bound for $\ell_2$

$S_1$ : a hyperspherical sector of a  $d$ -dimensional ball.

The classifier we defined is  $1_{\{x \in S_1\} \cup \{x \in S_2'\}}$  and  $\delta$  is orthogonal to the hyperplane.  $R_x$  is a sample and distribution based radius.

**2<sup>nd</sup> upper bound:**  $r < \|\delta\|_2 < \frac{5}{\sqrt{d}} R_x$



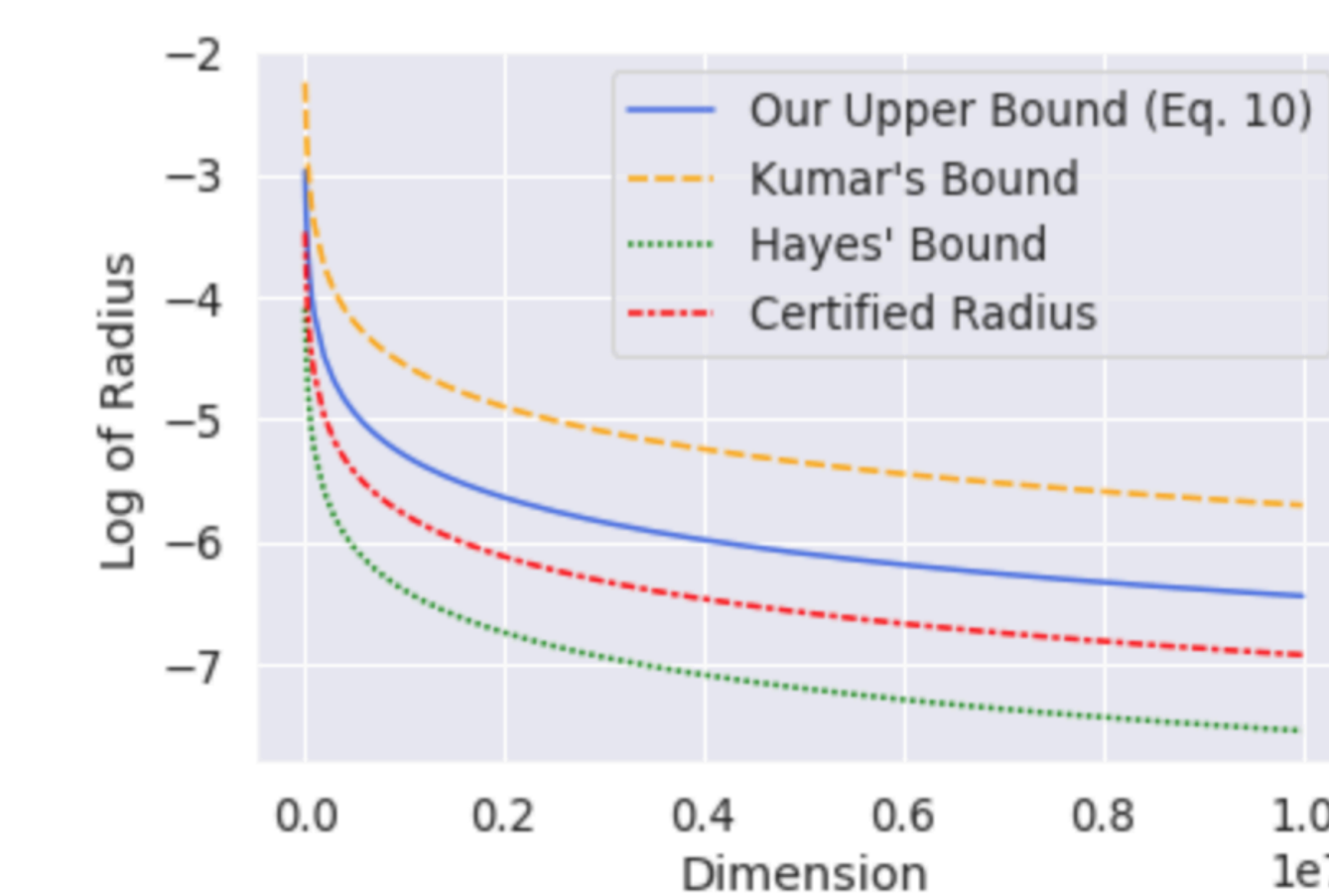
## Extension to $\ell_p$ bounds

We extend our  $\ell_2$  bounds to  $\ell_p$  bounds with spherical symmetric distributions.

$$r < \frac{5}{d^{1-1/p}} \Psi^{-1}\left(\frac{g(x)}{1 - 5 \times 10^{-7}}; q\right) \quad r < \frac{5}{d^{1-1/p}} R_x$$

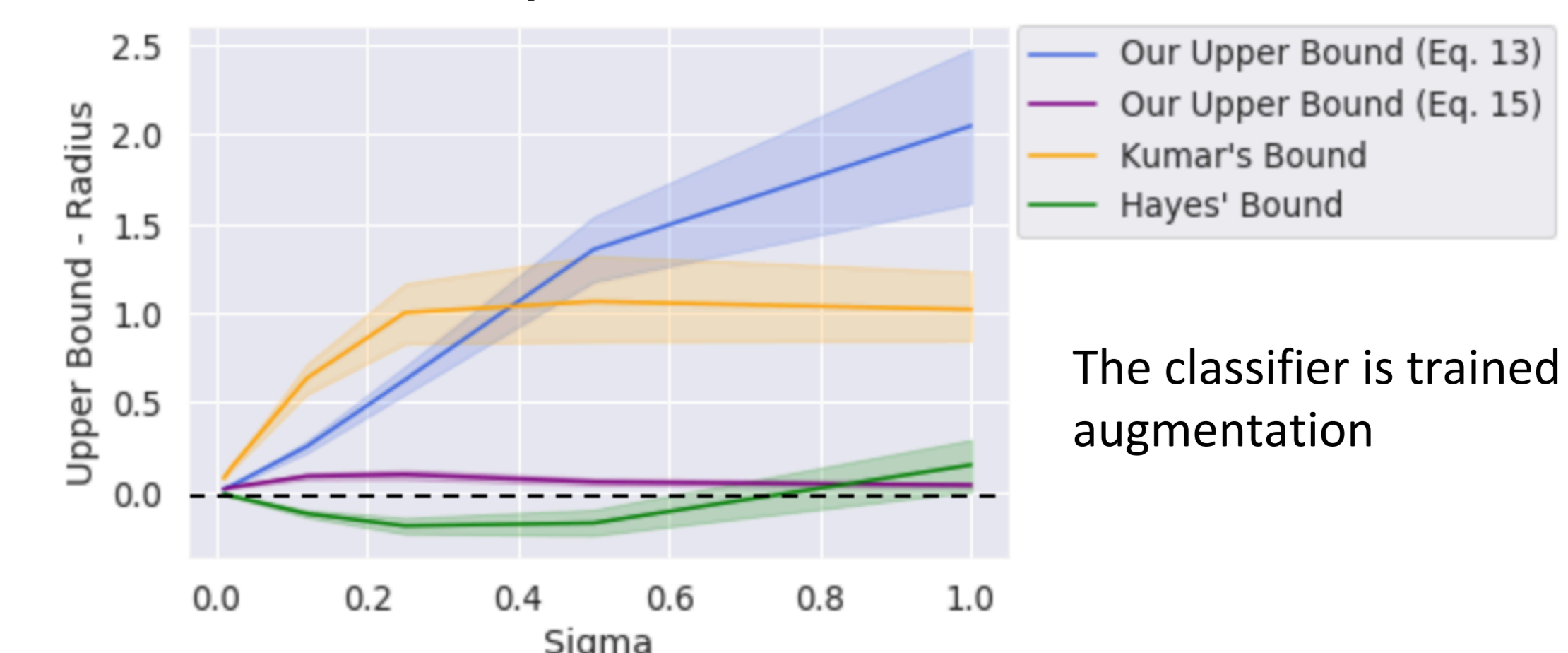
## Experiment with Gaussian smoothing

Evaluating our first upper bound for different parameters, we compare our result with SOTA bounds.



$$\sigma = 1/\sqrt{d}, g(x) = 0.999$$

Evaluating our second upper bound on CIFAR10 dataset with 100 random samples.



The classifier is trained with noise augmentation