# CS 6301.007

# Machine Learning in Cyber Security

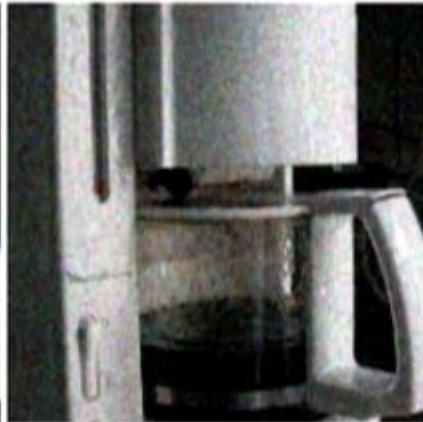Wei Yang

Department of Computer Science

University of Texas at Dallas

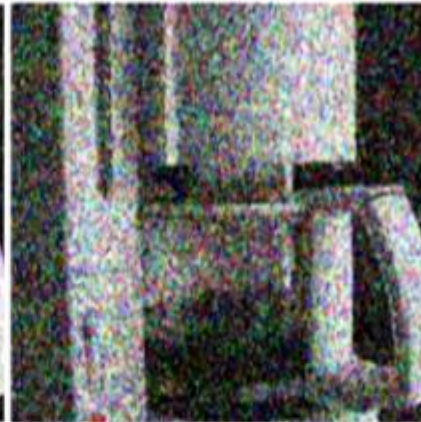What is the difference between adversarial examples and test error?

- InceptionV3: 13.2% test accuracy at sigma=.4
- (76.2% clean test accuracy)

# Salt-and-pepper noise



- InceptionV3: 5.4% test accuracy at p=.3

# Video Robustness



"Bus"
frame 3 pred 4 with p 0.27798435092

"Plane"
frame 3 pred 13 with p 0.823262214661

"Plane"
frame 3 pred 13 with p 0.914280056953
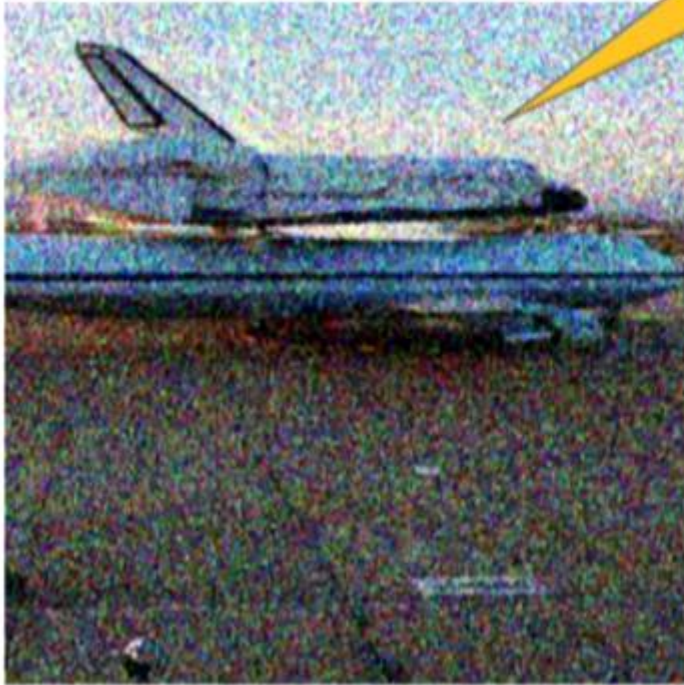
"Person"
frame 12 pred 0 with p 0.185599714518

Thanks to Keren Gu for this example!

# Naïve data augmentation doesn't help

# Compress noice != white noise



Original

Compressed

# Data augmentation can even hurt you
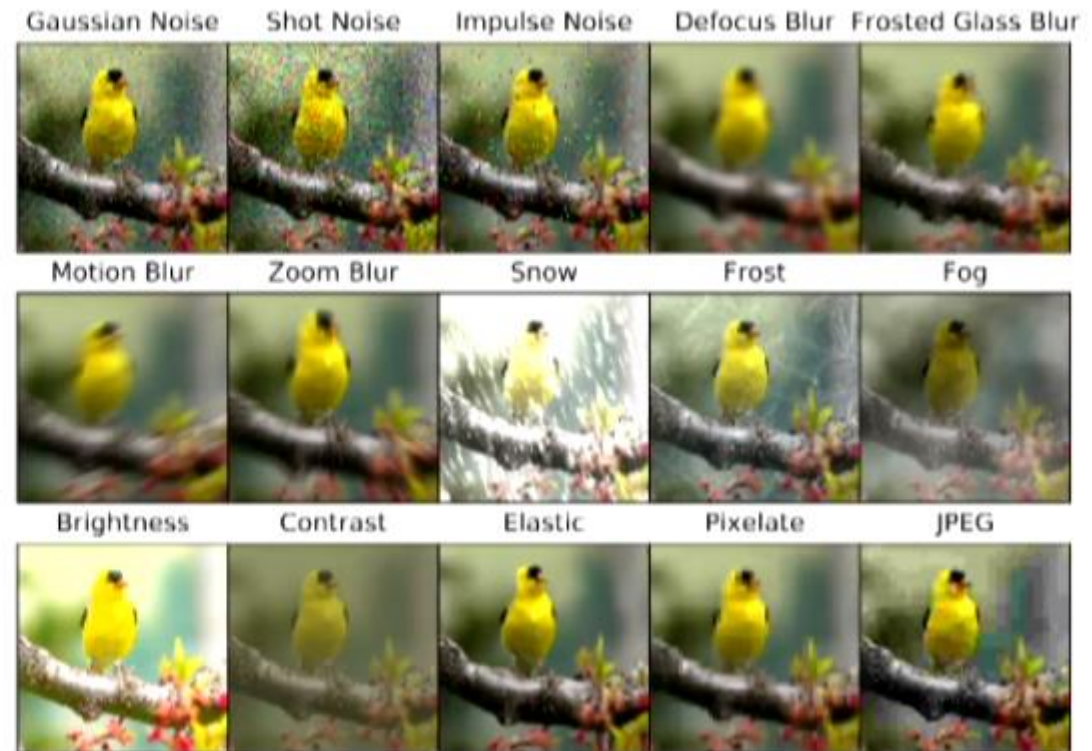
# Corruption Robustness

- Goal: Measure and improve model robustness to distributional shift.
- Corruptions are not worst-case.
- Test examples are randomly sampled to best estimate probability of an error.

Training on randomized textures helps

What is the difference between adversarial examples and test error?

# Adversarial Examples - Security



https://qz.com/721615/smart-pirates-are-fooling-youtubes-copyright-bots-by-hiding-movies-in-360-degree-videos/

# Questions for Design a Secure ML System

- How do adversaries typically break systems?
- How would you measure test error?
- Are you secure if test error > 0?
- How do we deal with out-of-distribution generalization?

# Adversarial Examples – ML Phenomenon



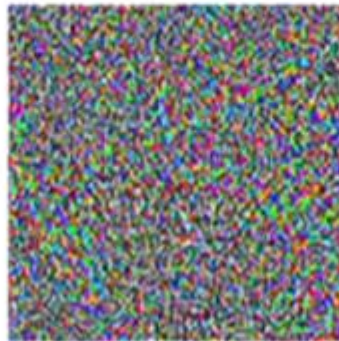**Why** do our models have adversarial examples?   **A:** ???

**What** are adversarial examples?   **A:** The nearest test error

"panda"
57.7% confidence

+ ε

=

"gibbon"
99.3% confidence

**Why** do our models have **test error?**

A: ???

**What** are adversarial examples?

A: The nearest test error



"panda"

57.7% confidence

$+\epsilon$

$=$

"gibbon"

99.3% confidence

**Why** do our models have **test error?**

A: ???

**What** are adversarial examples?

A: The nearest test error



Test error in what distribution?

"panda"
57.7% confidence

"gibbon"
99.3% confidence
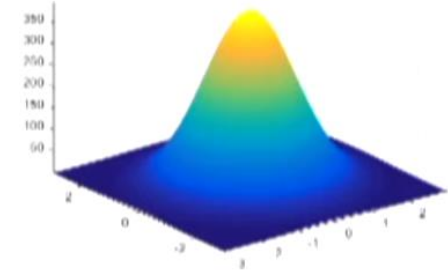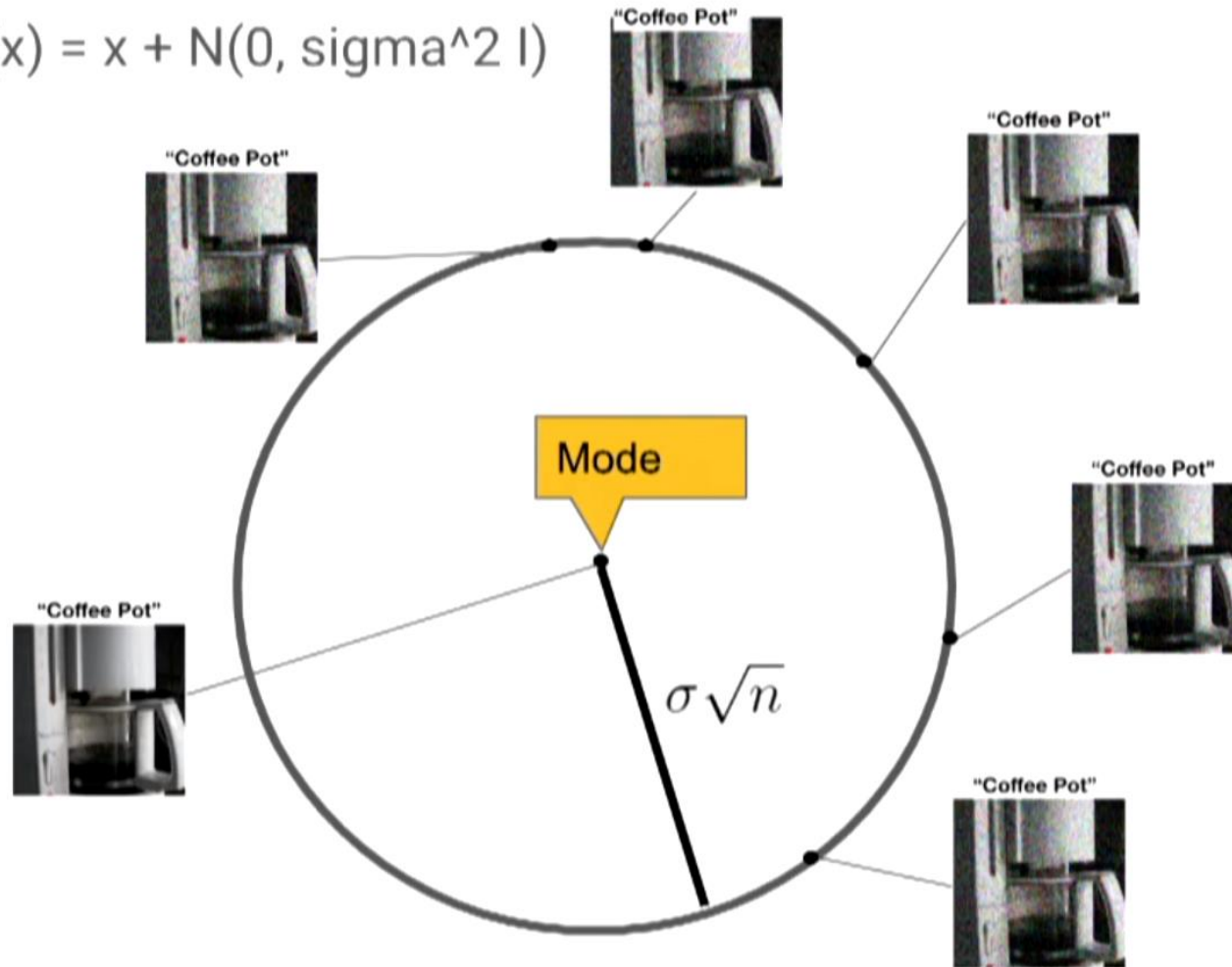
- In high dimensions, what does .1% test error look like?
- How close should the nearest test error be? (Assuming we sample infinitely)

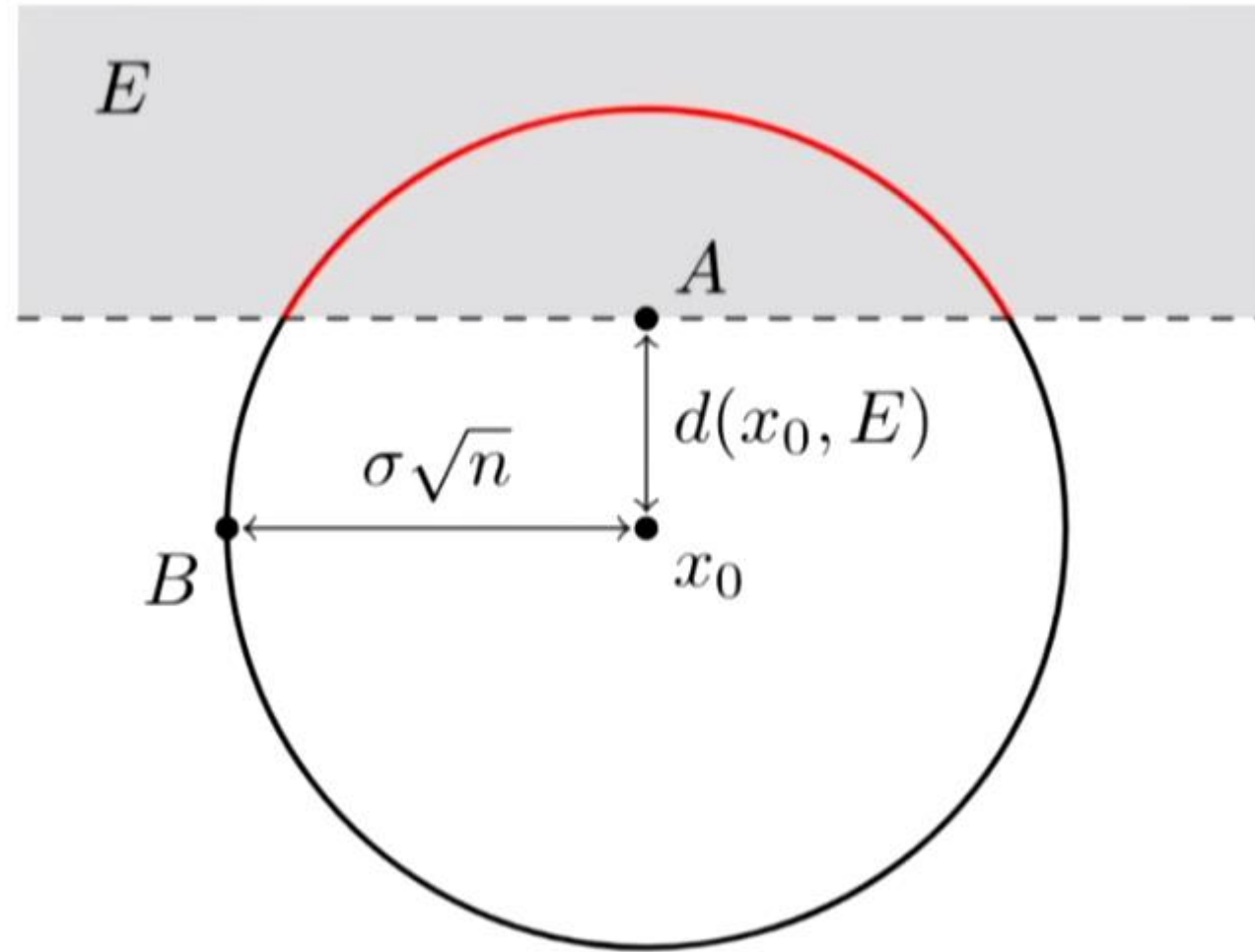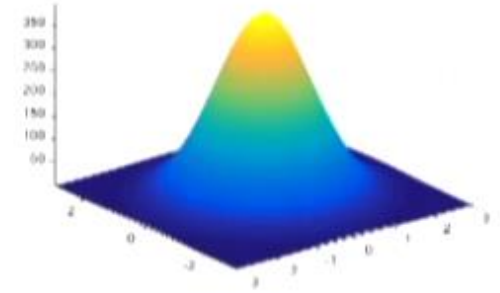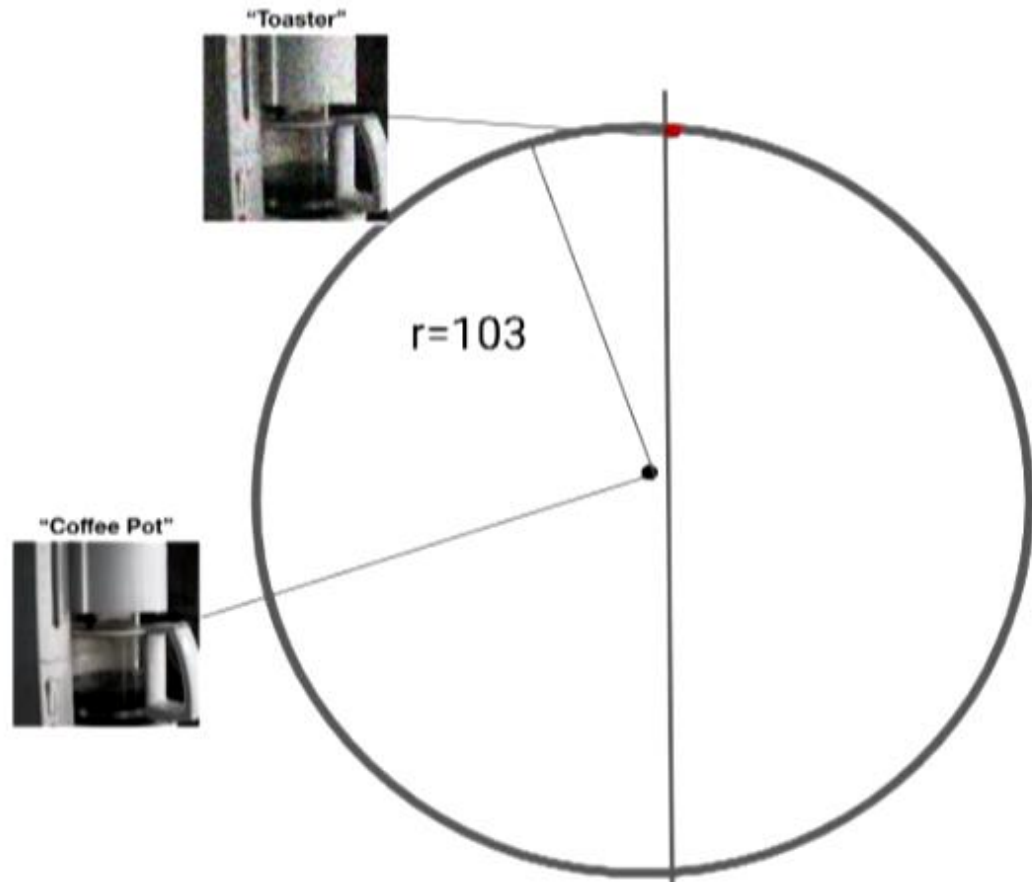# High Dimensional Gaussians



$q(x) = x + N(0, sigma^2 I)$

"Coffee Pot"

Mode

$\sigma \sqrt{n}$

- sigma=.2
- n=299*299*3
- 270,000 dimensional sphere
- radius ~ 103
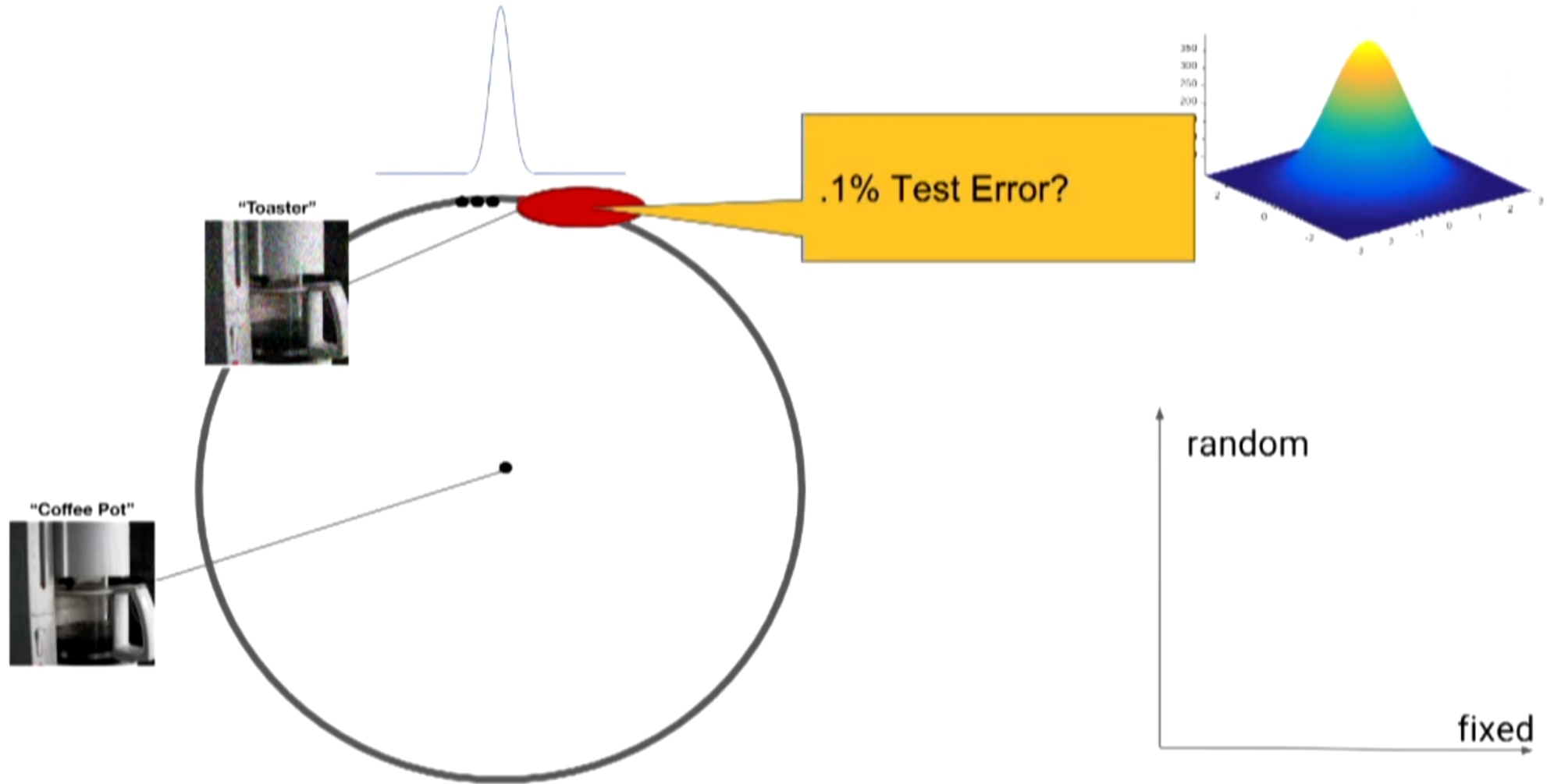
# Linear Models

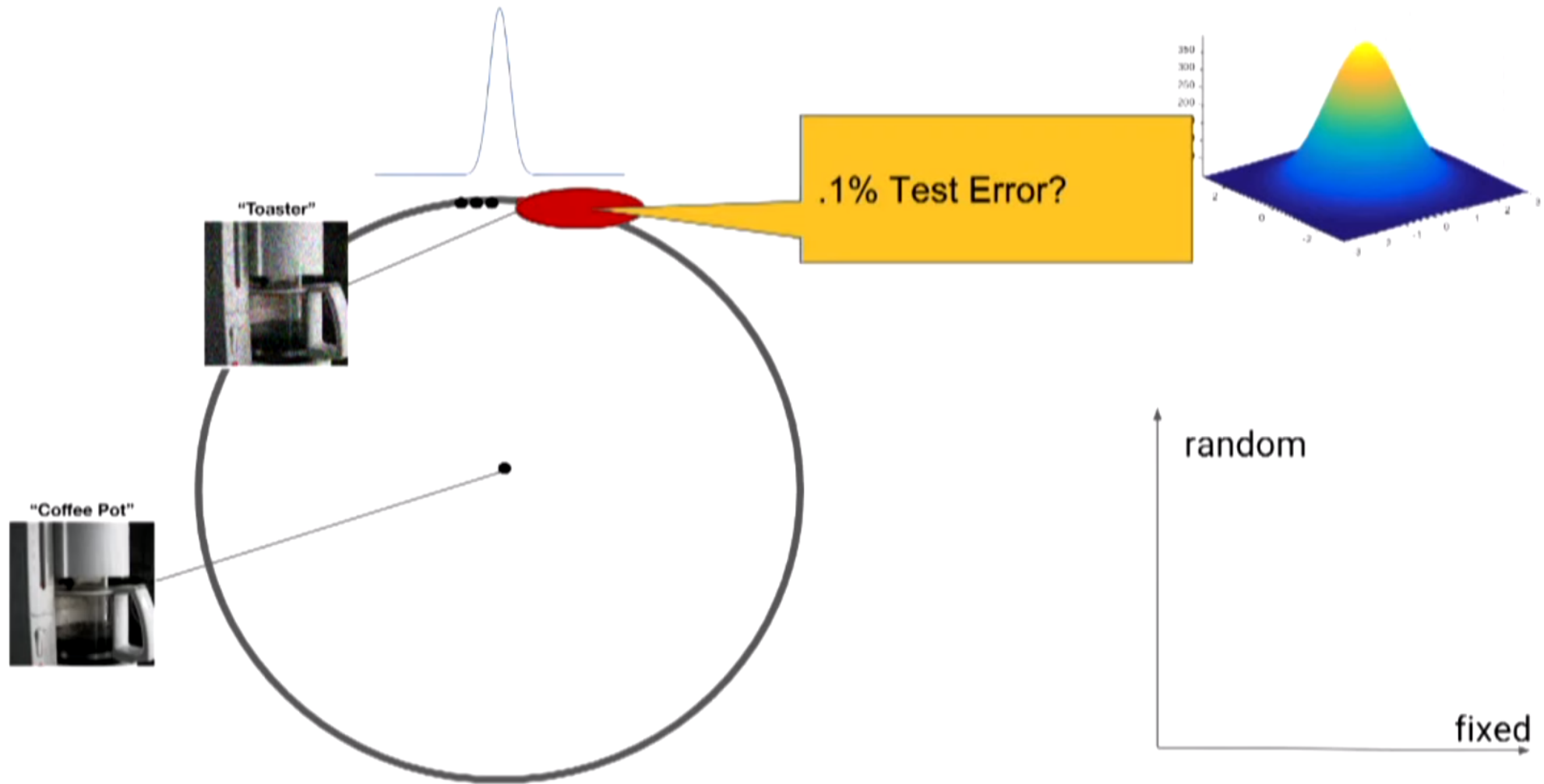$q(x) = x + N(0, sigma^2 I)$



"Toaster"

r=103

"Coffee Pot"



Theorem: A linear model with error rate mu in distribution q, has its nearest error at distance

$$\sigma\Phi^{-1}(\mu) = O(\sigma)$$

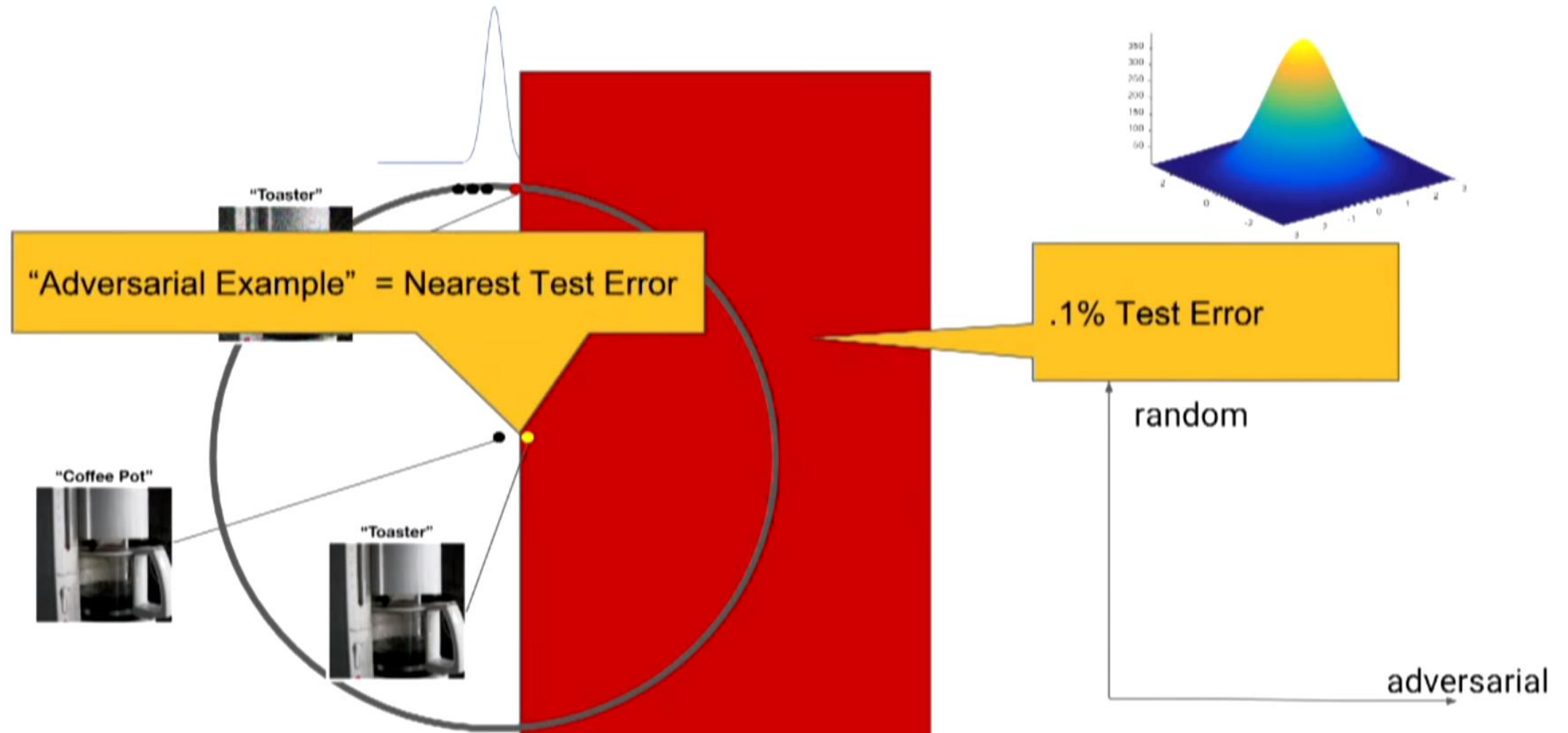- sigma=.2
- .1% error  ->  d = .62
- 10^-9 error  -> d = 1.2

https://arxiv.org/pdf/1608.08967.pdf

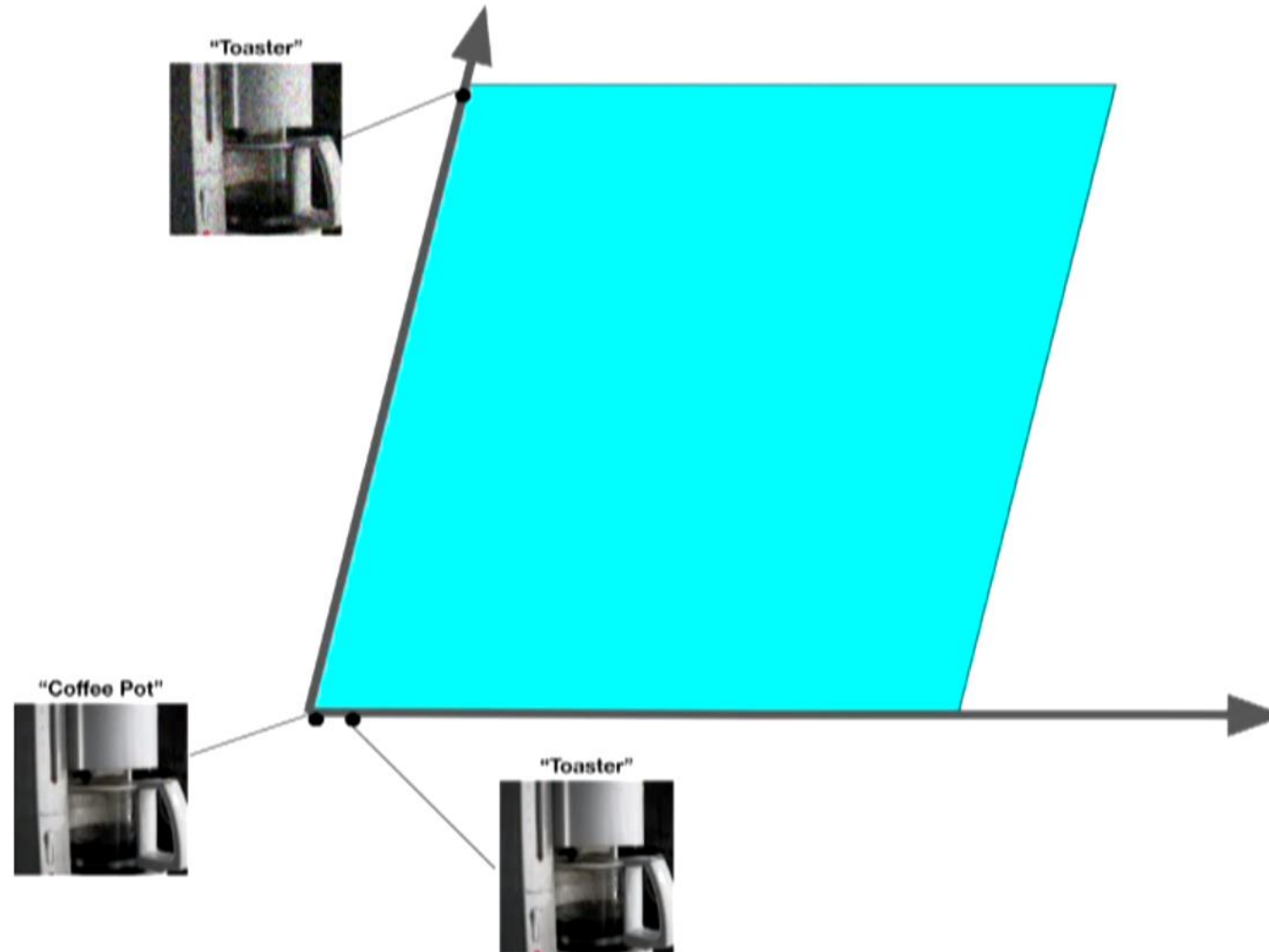# Where is .1% Test Error?

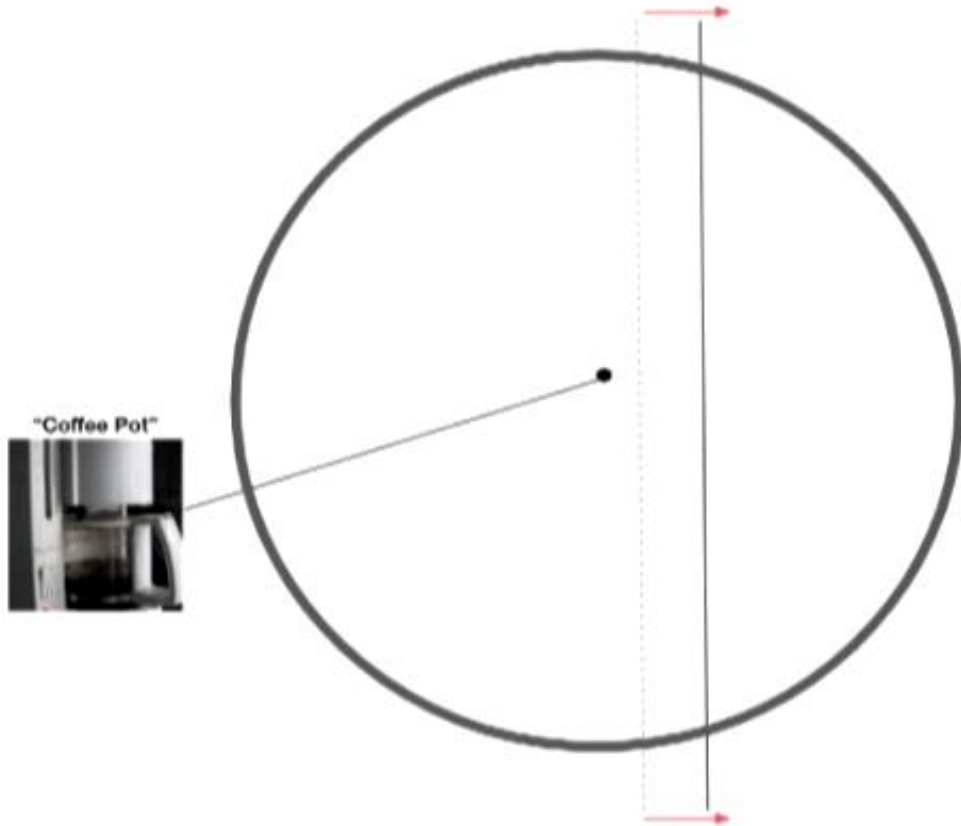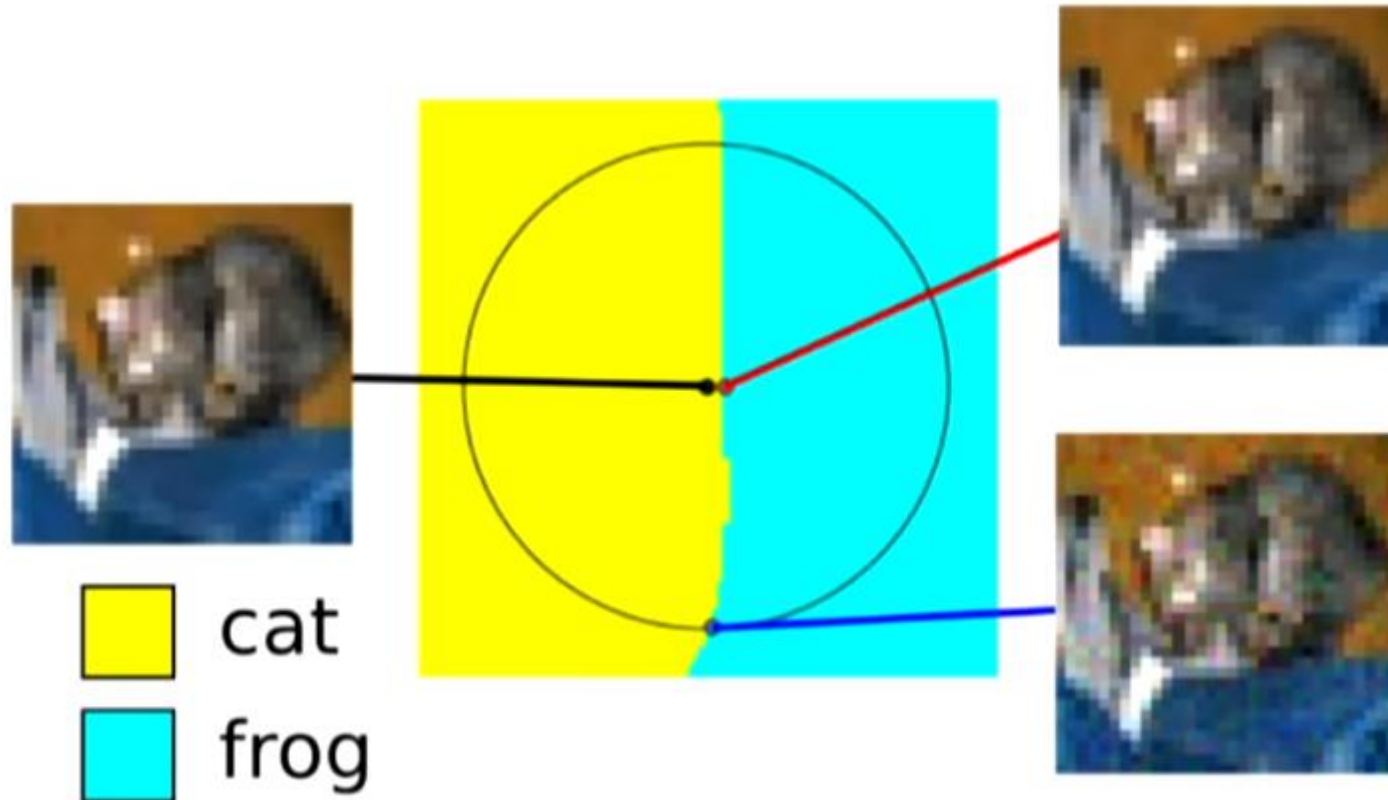# Where is .1% Test Error?

# Church Window Plot

# InceptionV3

σ vs. distance for clean points (ImageNet)

sigma=.04  (R=2.2)
error rate .2%
d = .16
predicted d = .08

# Resnet-50

"Cat"

σ vs. distance for clean points (CIFAR)

- linear
- naturally trained
- trained on noise

Distance to decision boundary

σ at which error rate is 0.01

# Resnet-50

"Cat"

σ vs. distance for clean points (CIFAR)

- linear
- naturally trained
- trained on noise
- adversarially trained

Distance to decision boundary

σ at which error rate is 0.01

Why are we trying to "defend" against the nearest test error?

"panda"
57.7% confidence

$$x_{adv} = \max_{x':\|x-x'\|_\infty < \epsilon} L(\theta, x', \hat{y})$$

Test error

L_inf eps=8/255

panda

miniature poodle

Tibetan mastiff

# Adversarial Defenses – Why?

# Successful Defenses Improve Robustness





https://arxiv.org/pdf/1706.06083.pdf

# Failed Defenses Don't Improve Robustness



Performance of broken adversarial defenses in noise

"[One] Possible explanation is that the **set of adversarial negatives** is of extremely low probability, and thus is never (or rarely) observed in the test set, yet it is dense **(much like the rational numbers)**, and so it is found near every virtually every test case."
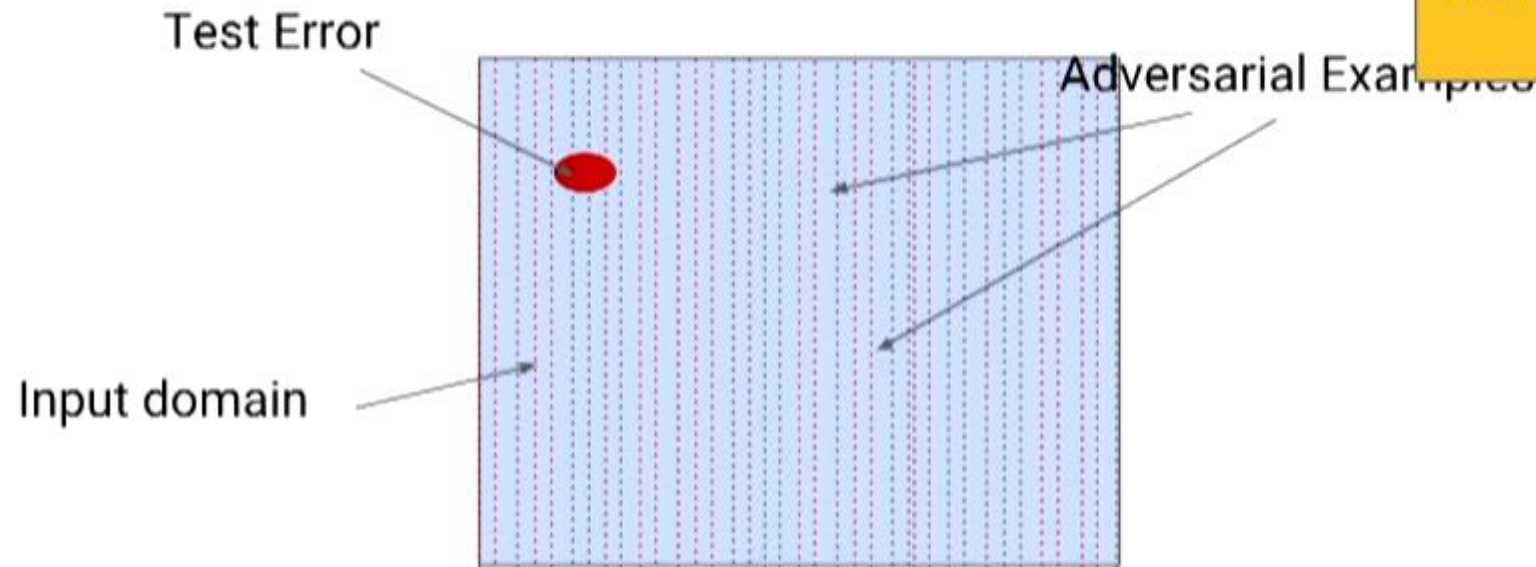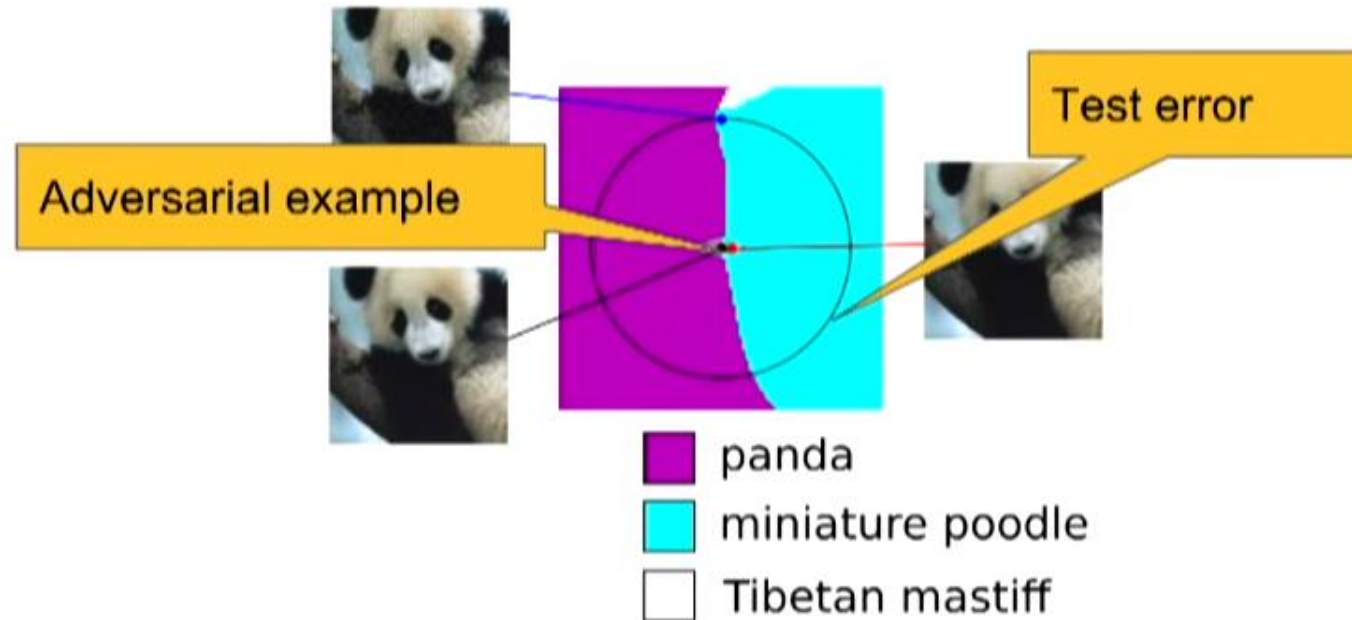
"[One] Possible explanation is that the **set of adversarial negatives** is of extremely low probability, and thus is never (or rarely) observed in the test set, yet it is dense **(much like the rational numbers)**, and so it is found near every virtually every test case."
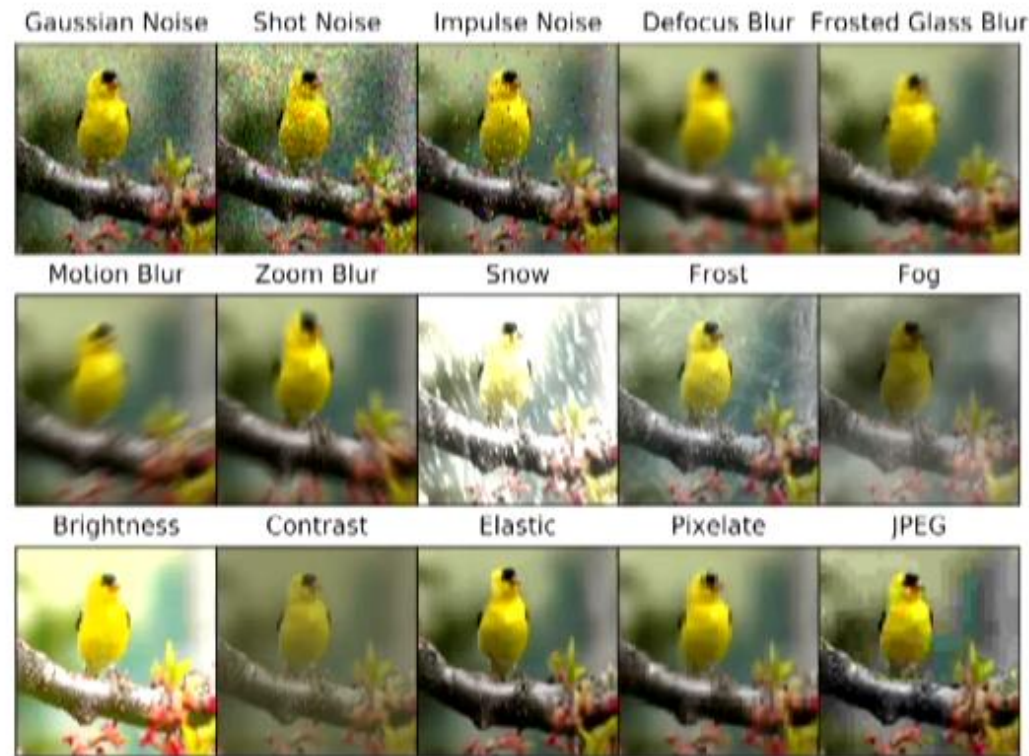
Make a harder test set

Test Error

Adversarial Examples

Input domain

- Adversarial examples are the nearest test error.
- Test error measures the **amount** of errors.
- The nearest error is not surprisingly close given the **amount** of errors.
- We can measure test error outside the natural distribution.
- **There is always going to be a nearest error.**

- Robustness to distributional shift is the *real* problem here.
- If you disagree, at least measure both for the sake of **science.**
- It's a critical sanity check for the vanishing gradient problem.



Hendrycks et. al.