

# Information-Theoretic Methods in Data Sciences

## Model Uncertainty, Robustness and Model Drift

Pablo Piantanida

pablo.piantanida@centralesupelec.fr

Laboratoire des Signaux et Systèmes (L2S)

Université Paris-Saclay CNRS CentraleSupélec

7th workshop “Journée Statistique et Informatique à Paris Saclay”,  
January 26th, 2022



CentraleSupélec

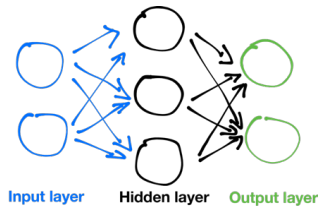
université  
PARIS-SACLAY



- 1 A Brief Overview of AI and Information Theory
  - The birth of AI and Deep Learning
  - Legacy of Shannon's work
  - Information, Uncertainty and Learning
- 2 Critical Problems in Safety AI
- 3 Overview of Recent Contributions to Safety AI
  - Detecting Misclassification Errors
  - Out-of-Distribution Detection
  - Adversarial Robustness
- 4 Discussion and Research Perspectives

- 1 A Brief Overview of AI and Information Theory
  - The birth of AI and Deep Learning
  - Legacy of Shannon's work
  - Information, Uncertainty and Learning
- 2 Critical Problems in Safety AI
- 3 Overview of Recent Contributions to Safety AI
  - Detecting Misclassification Errors
  - Out-of-Distribution Detection
  - Adversarial Robustness
- 4 Discussion and Research Perspectives

# The birth of AI and Deep Learning





# A Brief History of AI - Dartmouth Conference (1956)

*We propose that a 2 month, 10 man study of artificial intelligence be carried out during the summer of 1956 at Dartmouth College in Hanover, New Hampshire.*

*The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves.*

- [Dartmouth AI Project Proposal](#); J. McCarthy et al.; Aug. 31, 1955.



# A Brief History of AI - Dartmouth Conference (1956)



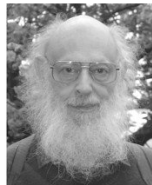
**John McCarthy**



**Marvin Minsky**



**Claude Shannon**



**Ray Solomonoff**



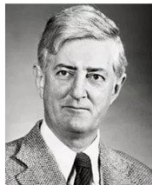
**Alan Newell**



**Herbert Simon**



**Arthur Samuel**



**Oliver Selfridge**



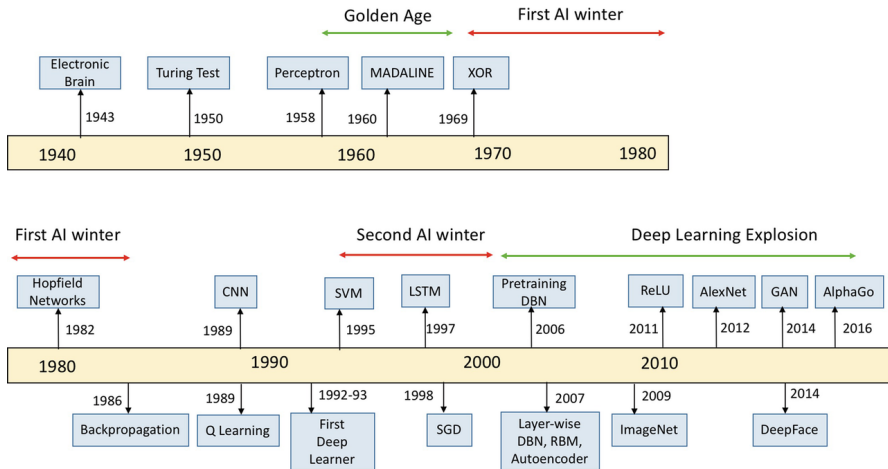
**Nathaniel Rochester**



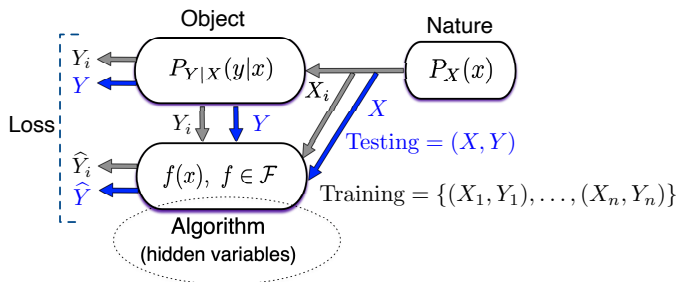
**Trenchard More**

## The Founding Fathers of AI

# A Brief History of AI - Timeline (1943 - Present)



# A Probabilistic Model of Learning (1960)



- **Imitation of the object:** try to construct a predictor which provides the best predictions to the supervisor output
- **Approximation of the object:** try to approximate the object (nature) itself based on a model (uncertainty and calibration)

**Learning is data compression:** To separate structure from noise, the regularities present in the data by choosing appropriately  $f \in \mathcal{F}$ .

$$P \left( \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)| > \varepsilon \right) \leq 8S(\mathcal{F}, n)e^{-n\varepsilon^2/32}$$
$$\mathbb{E} \left[ \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)| \right] \leq 2\sqrt{\frac{\log S(\mathcal{F}, n) + \log 2}{n}}$$



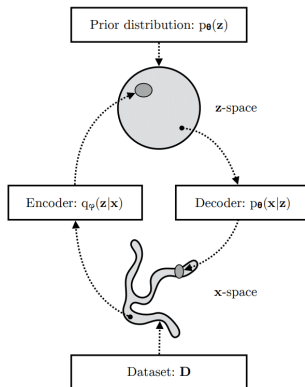
Vapnik–Chervonenkis theory (1960) addresses key questions:

- What are the conditions for **consistency of a learning rule** based on the empirical risk minimization principle?
- How fast is the **rate of convergence** of the learning process?
- How can one control the **generalization ability** (convergence rate) of the learning process?

Vapnik and Chervonenkis' ingenious formulation led to the characterization of **necessary and sufficient conditions** (finite VC-dimension) for the minimizing of a risk  $R(f)$  using data.

## Good representations learn to **disentangle manifolds**:

- Enc/Dec map between low and high representations of data,
- Encoders **perform inference** to **interpret data, flatten** and to **disentangle** the data manifold,
- Decoders can introduce changes in reconstructing data features,
- How **goodness** should be defined is an **open problem**.



My research focus on **developing and bringing new mathematical tools and methodological principles from information theory** to machine learning and deep learning.

## Legacy of Shannon's work

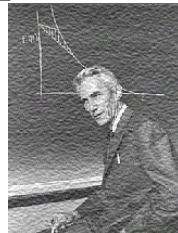


# Shannon's Model of a Communication System



Shannon proposed (1948) an asymptotic approach:

- A  $k$ -symbol sequence  $U$  is mapped by an encoder into an  $n$ -symbol input sequence  $X$
- The received channel output sequence  $Y$  is mapped by a decoder into an estimate  $\hat{U}$
- What is the the **maximum communication rate**  $R = k/n$  (bits per transmission) such that  $P\{U \neq \hat{U}\}$  **can be made arbitrarily small** when  $(k, n)$  are sufficiently large?



Shannon's ingenious formulation led to the characterization of **necessary and sufficient conditions** for reliable communication.





Shannon proposed (1948) an  $(k, n)$ -asymptotic approach:

- **Channel coding theorem:** The capacity is the maximum of the **mutual information** between the channel input and output

$$C = \sup_{p_X} I(X; Y) \quad \text{in bit/transmission.}$$

- **Lossy source coding theorem:** The optimal tradeoff between the rate  $R = n/k$  and the distortion  $D$  is

$$R(D) = \inf_{p_{\hat{U}|U}: \mathbb{E}[d(\hat{U}, U)] \leq D} I(U; \hat{U}) \quad \text{in bits/symbol.}$$

- **Separation theorem:** Shannon's ingenious formulation led to

$$R(D) < C,$$

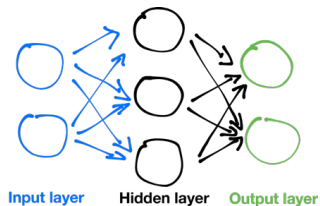
**necessary and sufficient conditions** for reliable communication.

*“Using bits as a universal representation between sources and channels is essentially optimal”*

Nothing is more practical than a good theory:

- Analogue **data can be represented by discrete symbols** and compressed before transmission
- **Representation of information is at the heart of modern communications** (codes that can squish messages, saving time resources and codes that can protect data from noise)
- Information theory provides valuable insight, highlighting key properties of good codes, **leading to optimal designs.**

# Information, Uncertainty and Learning



**Entropy**  $H(X)$  of a discrete random variable (RV)  $X \sim p$ :

- **1. Measure of uncertainty**  $\rightarrow$  “surprise” function  $s(x)$ ,  $x \in \mathcal{X}$ , and  $H(X) = \mathbb{E}[s(X)]$
- **2. Independent of alphabet**  $\rightarrow s(x) = s(p(x))$
- **3. Additivity:**

$$s(p(x)q(y)) = s(p(x)) + s(q(y)) \quad \rightarrow \quad s(x) = \log p(x)$$

- Lower probability implies higher surprise  $\rightarrow s(x) = -\log p(x)$

$$\begin{aligned} H(X) &= - \sum_{x \in \mathcal{X}} p(x) \log p(x) \\ &= -\mathbb{E}[\log p(X)] \end{aligned}$$

- $H(X)$  is nonnegative, continuous, and strictly concave function of  $p$ , and  $0 \leq H(X) \leq \log |\mathcal{X}|$ .

**Entropy**  $H(X)$  of a discrete random variable (RV)  $X \sim p$ :

- **1. Measure of uncertainty**  $\rightarrow$  “surprise” function  $s(x)$ ,  $x \in \mathcal{X}$ , and  $H(X) = \mathbb{E}[s(X)]$
- **2. Independent of alphabet**  $\rightarrow s(x) = s(p(x))$
- **3. Additivity:**

$$s(p(x)q(y)) = s(p(x)) + s(q(y)) \quad \rightarrow \quad s(x) = \log p(x)$$

- Lower probability implies higher surprise  $\rightarrow s(x) = -\log p(x)$

$$\begin{aligned} H(X) &= - \sum_{x \in \mathcal{X}} p(x) \log p(x) \\ &= -\mathbb{E}[\log p(X)] \end{aligned}$$

- $H(X)$  is nonnegative, continuous, and strictly concave function of  $p$ , and  $0 \leq H(X) \leq \log |\mathcal{X}|$ .

- Rényi entropy for a discrete r.v.  $X$  with probability  $p(x)$ :

$$\begin{aligned}H_{\alpha}(X) &= \frac{1}{1-\alpha} \log \sum_{x \in \mathcal{X}} p(x)^{\alpha} \\ &= \frac{1}{1-\alpha} \log \mathbb{E}[p(X)^{\alpha-1}],\end{aligned}$$

for  $\alpha > 0$ ;  $H_{\alpha}(X) \rightarrow H(X)$  as  $\alpha \rightarrow 1$ .

- Conditional Rényi entropy for discrete RVs  $(X, Y) \sim p(x, y)$ :

$$\begin{aligned}H_{\alpha}(X|Y) &= \sum_{y \in \mathcal{Y}} p(y) \left( \frac{1}{1-\alpha} \log \sum_{x \in \mathcal{X}} p(x|y)^{\alpha} \right) \\ &= \frac{1}{1-\alpha} \mathbb{E} \left[ \log \sum_{x \in \mathcal{X}} p(x|Y)^{\alpha} \right].\end{aligned}$$

- There are many other information measures.

- Rényi entropy for a discrete r.v.  $X$  with probability  $p(x)$ :

$$\begin{aligned}H_{\alpha}(X) &= \frac{1}{1-\alpha} \log \sum_{x \in \mathcal{X}} p(x)^{\alpha} \\ &= \frac{1}{1-\alpha} \log \mathbb{E}[p(X)^{\alpha-1}],\end{aligned}$$

for  $\alpha > 0$ ;  $H_{\alpha}(X) \rightarrow H(X)$  as  $\alpha \rightarrow 1$ .

- Conditional Rényi entropy for discrete RVs  $(X, Y) \sim p(x, y)$ :

$$\begin{aligned}H_{\alpha}(X|Y) &= \sum_{y \in \mathcal{Y}} p(y) \left( \frac{1}{1-\alpha} \log \sum_{x \in \mathcal{X}} p(x|y)^{\alpha} \right) \\ &= \frac{1}{1-\alpha} \mathbb{E} \left[ \log \sum_{x \in \mathcal{X}} p(x|Y)^{\alpha} \right].\end{aligned}$$

- There are many other information measures.

## An emerging interface?

- Shannon's entropy provides a **measure of uncertainty** about the amount of information that a learner possesses **relative to a given concept** when only the probability distribution is given.
- **But** the basic problem of learning consists in that **one has to separate the relevant information** from patterns.

## Two questions naturally arise:

- Are information measures **fundamental measures of the random properties of data** for learning problems?
- What are the instances of learning problems for which **information measures can play a key role**?

The study of these questions has played an important role and, undoubtedly, it will play a central role in future learning methods.



# A Long-Lasting Partnership

Learning is data compression:

- The goal is to learn the laws and regularities present in the data, that is, **to separate structure from noise**.
- **Data compression is fundamentally related to statistical generalization** as shown by a number of sample complexity bounds (e.g., VC-dimension, PAC-Bayes, and others). 😊
- The celebrated **Minimum Description Length (MDL)** principle, to approach model selection in statistical inference.
- In unsupervised learning, **Variational Autoencoders (VAEs)** are motivated by compression methods and the **Information Bottleneck** method for supervised learning as well.

This talk focuses on two related problems:

- **How to measure uncertainty from model predictions?**
- **How to detect uncertainty induced from data drift?**

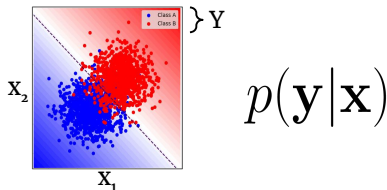
- 1 A Brief Overview of AI and Information Theory
  - The birth of AI and Deep Learning
  - Legacy of Shannon's work
  - Information, Uncertainty and Learning
- 2 Critical Problems in Safety AI
- 3 Overview of Recent Contributions to Safety AI
  - Detecting Misclassification Errors
  - Out-of-Distribution Detection
  - Adversarial Robustness
- 4 Discussion and Research Perspectives

# What Does Model Uncertainty Means?

Return a distribution over predictions rather than a single prediction.

- **Classification:** Output label along with its confidence.
- **Regression:** Output mean along with its variance.

Good uncertainty estimates quantify *when we can trust the model's predictions*.



$$p(\mathbf{y}|\mathbf{x})$$

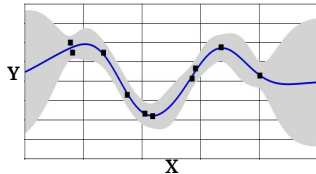


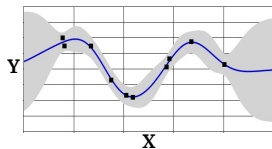
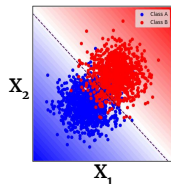
Image credit: Eric Nalisnick

# What Does Out-of-Distribution Robustness Means?

**I.I.D.**  $p_{\text{TEST}}(y,x) = p_{\text{TRAIN}}(y,x)$

*(Independent and Identically Distributed)*

**O.O.D.**  $p_{\text{TEST}}(y,x) \neq p_{\text{TRAIN}}(y,x)$



*Image credit: Eric Nalisnick*

# What Does Out-of-Distribution Robustness Means?

**I.I.D.**       $p_{\text{TEST}}(y,x) = p_{\text{TRAIN}}(y,x)$

**O.O.D.**       $p_{\text{TEST}}(y,x) \neq p_{\text{TRAIN}}(y,x)$

Examples of dataset shift:

- **Covariate shift.** Distribution of features  $p(x)$  changes and  $p(y|x)$  is fixed.
- **Open-set recognition.** New classes may appear at test time.
- **Label shift.** Distribution of labels  $p(y)$  changes and  $p(x|y)$  is fixed.

# ImageNet-C: Varying Intensity for Dataset Shift

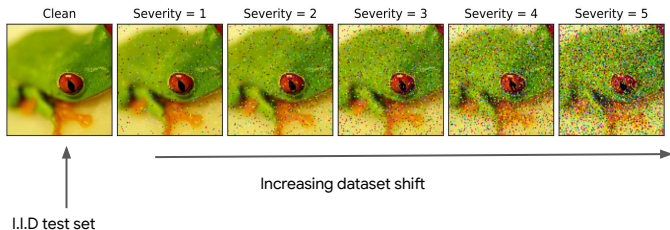


Image source: Benchmarking Neural Network Robustness to Common Corruptions and Perturbations, [Hendrycks & Dietterich, 2019](#).

# ImageNet-C: Varying Intensity for Dataset Shift

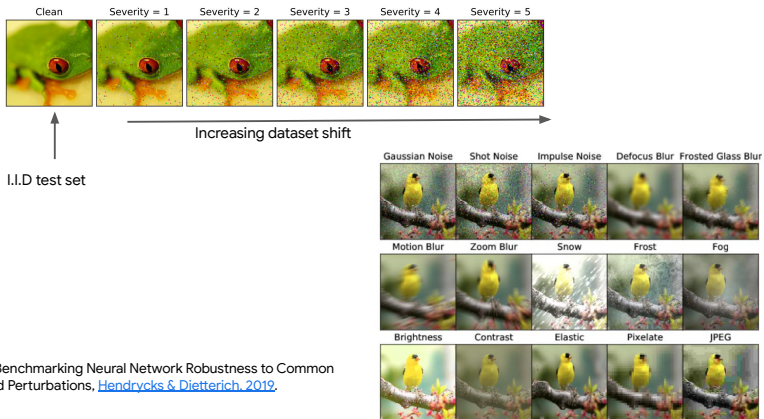
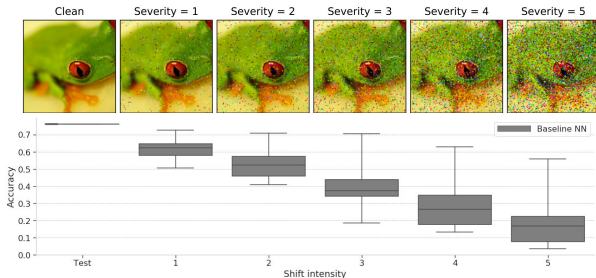


Image source: Benchmarking Neural Network Robustness to Common Corruptions and Perturbations, [Hendrycks & Dietterich, 2019](#).

# Neural Networks Do Not Generalize Under Distribution Shift

- **Accuracy drops** with increasing shift on Imagenet-C



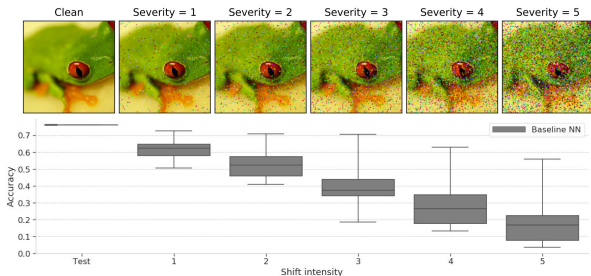
- But do the models know that they are less accurate?

Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift?, [Ovadia et al. 2019](#)

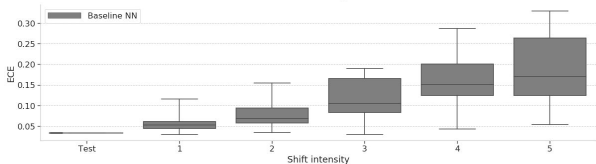


# Neural Networks Do Not Know When They Are Wrong

- **Accuracy drops** with increasing shift on Imagenet-C



- **Quality of uncertainty degrades** with shift  
-> “overconfident mistakes”



# Models Assign High Confidence Predictions to OOD Inputs

Example images where model assigns >99.5% confidence.

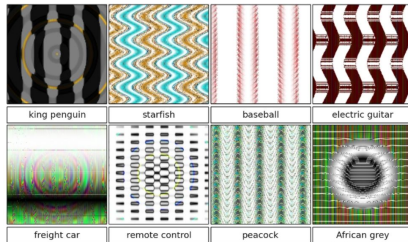
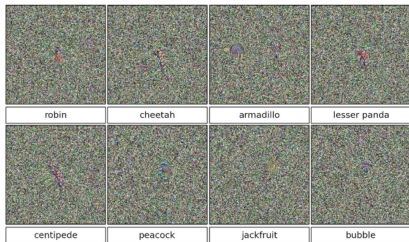


Image source: "Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images" [Nguyen et al. 2014](#)

# Models Assign High Confidence Predictions to OOD Inputs

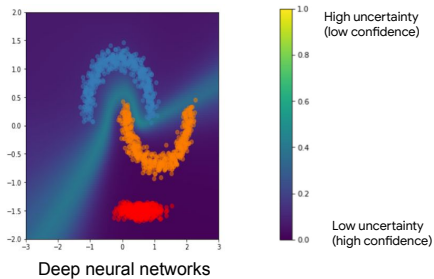
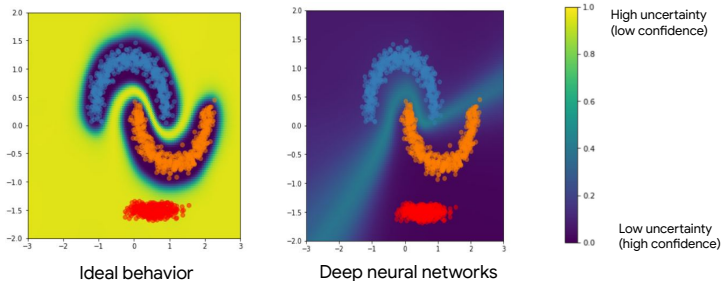


Image source: "Simple and Principled Uncertainty Estimation with Deterministic Deep Learning via Distance Awareness" [Liu et al. 2020](#)

# Models Assign High Confidence Predictions to OOD Inputs

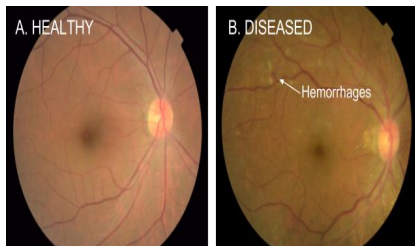


Trust model when  $x^*$  is close to  $p_{\text{TRAIN}}(x,y)$

Image source: "Simple and Principled Uncertainty Estimation with Deterministic Deep Learning via Distance Awareness" [Liu et al. 2020](#)

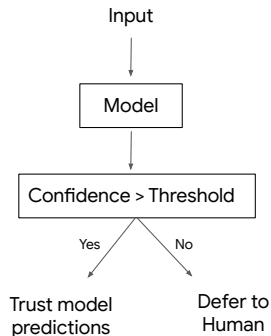
# Applications to Healthcare

- Use model uncertainty to decide when to trust the model or to defer to a human.
- Reject low-quality inputs.

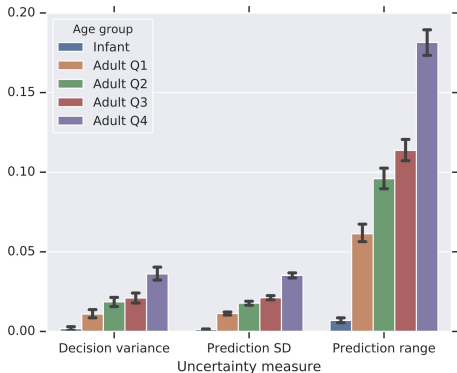


Diabetic retinopathy detection from fundus images

[Gulshan et al, 2016](#)



- Model accuracy and uncertainty across patient sub-groups



Mortality prediction from electronic health records

[Dusenberry et al. 2020](#)

# Applications to Self-driving Cars

Dataset shift:

- Time of day / Lighting
- Geographical location (City vs suburban)
- Changing conditions (Weather / Construction)



Weather



Construction

Image credit: Sun et al. [Waymo Open Dataset](#)



Daylight



Night



Downtown



Suburban

# Applications to Open Set Recognition

- Example: Classification of genomic sequences

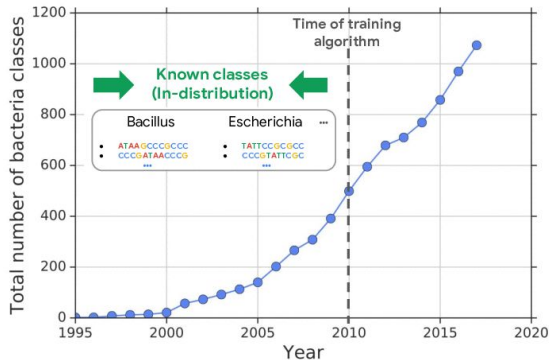


Image source: <https://ai.googleblog.com/2019/12/improving-out-of-distribution-detection.html>



- 1 A Brief Overview of AI and Information Theory
  - The birth of AI and Deep Learning
  - Legacy of Shannon's work
  - Information, Uncertainty and Learning
- 2 Critical Problems in Safety AI
- 3 Overview of Recent Contributions to Safety AI
  - Detecting Misclassification Errors
  - Out-of-Distribution Detection
  - Adversarial Robustness
- 4 Discussion and Research Perspectives

# DOCTOR: A Simple Method for Detecting Misclassification Errors

Joint work with Federica Granese, Marco Romanelli,  
Daniele Gorla and Catuscia Palamidessi



<https://neurips.cc/virtual/2021/spotlight/28017>

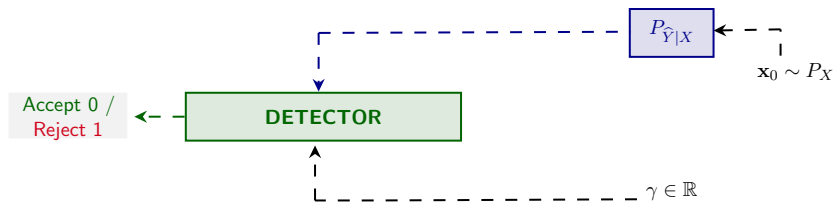
# Main Definitions

Let

- \*  $\mathcal{X} \subseteq \mathbb{R}$  be the **feature space**;
- \*  $\mathcal{Y} = \{1, \dots, C\}$  be the **label space**;
- \*  $p_{XY}$  be the underlying (unknown) probability density function over  $\mathcal{X} \times \mathcal{Y}$ ;
- \*  $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim p_{XY}$  be a random realization of  $n$  i.i.d. samples according to  $p_{XY}$  denoting the **training set**;
- \*  $f_{\mathcal{D}_n} : \mathcal{X} \rightarrow \mathcal{Y}$  be the **predictor**,

$$f_{\mathcal{D}_n}(\mathbf{x}) \equiv f_n(\mathbf{x}; \mathcal{D}_n) \triangleq \arg \max_{y \in \mathcal{Y}} P_{\hat{Y}|X}(y|\mathbf{x}; \mathcal{D}_n).$$

# Problem Definition

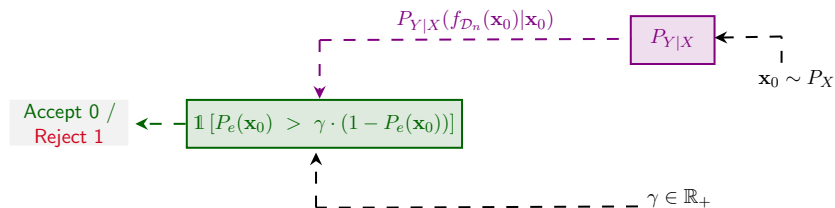


# Ideal (Oracle) Detector

## Definition (Error probability per sample)

For a given testing feature  $\mathbf{x}_0 \in \mathcal{X}$ ,

- \*  $E(\mathbf{x}_0) \triangleq \mathbb{1}[Y \neq f_{\mathcal{D}_n}(\mathbf{x}_0)]$  is the **error variable** corresponding to a predetermined predictor  $f_{\mathcal{D}_n}$  (based on  $P_{Y|X}$ );
- \*  $P_e(\mathbf{x}_0) \triangleq \mathbb{E}[E(\mathbf{x}_0)|\mathbf{x}_0] = 1 - P_{Y|X}(f_{\mathcal{D}_n}(\mathbf{x}_0)|\mathbf{x}_0)$  is the **probability of error classification w.r.t.  $P_{Y|X}$** .



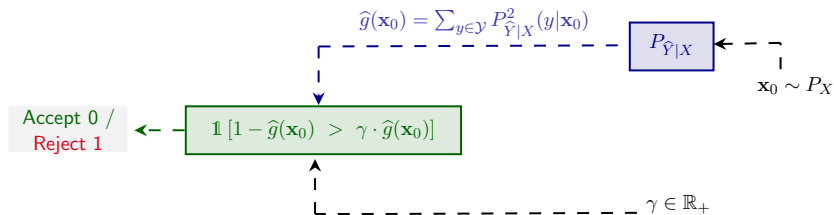
In practice,  $P_e(\mathbf{x}_0)$  **is not available**, but can we approximate it?

Proposition (DOCTOR:  $D_\alpha$ )

For a given testing feature  $\mathbf{x}_0 \in \mathcal{X}$ ,

- \*  $1 - \hat{g}(\mathbf{x}_0) \triangleq \sum_{y \in \mathcal{Y}} P_{\hat{Y}|X}(y|\mathbf{x}_0) \Pr(\hat{Y} \neq y|\mathbf{x}_0) = 1 - \sum_{y \in \mathcal{Y}} P_{\hat{Y}|X}^2(y|\mathbf{x}_0)$   
approximates the **probability of incorrectly classifying**  $\mathbf{x}_0$ ;
- \*  $(1 - \sqrt{\hat{g}(\mathbf{x}_0)}) - \Delta(\mathbf{x}_0) \leq P_e(\mathbf{x}_0) \leq (1 - \sqrt{\hat{g}(\mathbf{x}_0)}) + \Delta(\mathbf{x}_0)$  where

$$\Delta(\mathbf{x}_0) \triangleq 2\sqrt{2 \mathbf{KL}(P_{Y|X}(\cdot|\mathbf{x}_0) \| P_{\hat{Y}|X}(\cdot|\mathbf{x}_0))}.$$

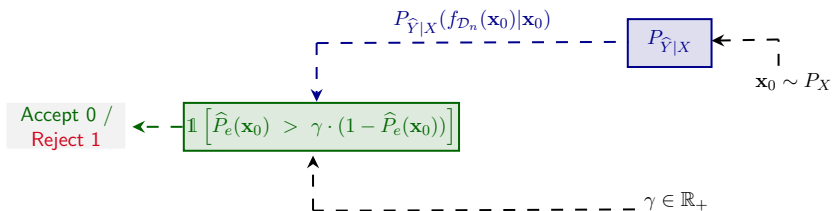


$$D_\alpha(\mathbf{x}_0, \gamma) \triangleq 1 [1 - \hat{g}(\mathbf{x}_0) > \gamma \cdot \hat{g}(\mathbf{x}_0)]$$

Definition (DOCTOR:  $D_\beta$ )

For a given testing feature  $\mathbf{x} \in \mathcal{X}$ ,

- \*  $\hat{E}(\mathbf{x}_0) \triangleq \mathbb{1}[\hat{Y} \neq f_{\mathcal{D}_n}(\mathbf{x}_0)]$  is the **self-error variable** corresponding to  $f_{\mathcal{D}_n}$  (based on the model  $P_{\hat{Y}|X}$ );
- \*  $\hat{P}_e(\mathbf{x}_0) \triangleq \mathbb{E}[\hat{E}(\mathbf{x}_0)|\mathbf{x}_0] = 1 - P_{\hat{Y}|X}(f_{\mathcal{D}_n}(\mathbf{x}_0)|\mathbf{x}_0)$  is the **probability of error classification w.r.t.  $P_{\hat{Y}|X}$** .



$$D_\beta(\mathbf{x}_0, \gamma) \triangleq \mathbb{1} \left[ \hat{P}_e(\mathbf{x}_0) > \gamma \cdot (1 - \hat{P}_e(\mathbf{x}_0)) \right]$$

## Definition (**FRR versus TRR**)

The false rejection rate (FRR) represents the probability that a hit (sample correctly classified) is rejected, while the true rejection rate (TRR) is the probability that a miss (sample wrongly classified) is rejected.

## Definition (**AUROC**)

The area under the Receiver Operating Characteristic curve (ROC) depicts the relationship between TRR and FRR. The perfect detector corresponds to a score of 100%.

## Definition (**FRR at 95% TRR**)

This is the probability that a hit is rejected when the TRR is at 95%.



# Scenarios: Totally Black Box & Partially Black Box

## Definition (**Totally Black Box (TBB) Scenario**)

In TBB only the output of the last layer of the network is available, hence gradient-propagation to perform input pre-processing is not allowed.

## Definition (**Partially Black Box (PBB) Scenario**)

In PBB we allow method-specific inputs perturbations and the possibility of doing temperature scaling.

## 1) ODIN [Liang et al., 2018]

$$\mathbf{SODIN}(\tilde{\mathbf{x}}) = \max_{i=[1:C]} \frac{\exp(f_i(\tilde{\mathbf{x}})/T)}{\sum_{j=1}^C \exp(f_j(\tilde{\mathbf{x}})/T)}$$
$$\mathbf{ODIN}(\tilde{\mathbf{x}}; \delta, T, \epsilon) = \begin{cases} \text{out}, & \text{if } \mathbf{SODIN}(\tilde{\mathbf{x}}) \leq \delta \\ \text{in}, & \text{if } \mathbf{SODIN}(\tilde{\mathbf{x}}) > \delta \end{cases}$$

- \*  $f(\tilde{\mathbf{x}})$  the vector of logits;
- \*  $\tilde{\mathbf{x}}$  represents a magnitude  $\epsilon$  perturbation of the original  $\mathbf{x}$ ;
- \*  $T$  is the temperature scaling parameter;
- \*  $\delta \in [0, 1]$  is the threshold value;
- \* *in* indicates the acceptance decision;
- \* *out* indicates the rejection decision.

## 2) Mahalanobis distance [Lee et al., 2018]

$$\mathbf{M}(\tilde{\mathbf{x}}) = \max_{c \in \mathcal{Y}} -(f(\tilde{\mathbf{x}}) - \hat{\mu}_c)^\top \hat{\Sigma}^{-1} (f(\tilde{\mathbf{x}}) - \hat{\mu}_c)$$
$$\text{MHLNB}(\tilde{\mathbf{x}}; \zeta, \epsilon) = \begin{cases} \text{out}, & \text{if } \mathbf{M}(\tilde{\mathbf{x}}) > \zeta \\ \text{in}, & \text{if } \mathbf{M}(\tilde{\mathbf{x}}) \leq \zeta \end{cases}$$

- ✱  $\hat{\mu}_c$  is the *empirical class mean* for each class  $c$  (training set);
- ✱  $\hat{\Sigma}$  is the *empirical covariance* (training set);
- ✱  $f(\tilde{\mathbf{x}})$  the vector of logits;
- ✱  $\tilde{\mathbf{x}}$  represents a magnitude  $\epsilon$  perturbation of the original  $\mathbf{x}$ ;
- ✱  $\zeta \in \mathbb{R}_+$  is the threshold value;
- ✱ *in* indicates the acceptance decision;
- ✱ *out* indicates the rejection decision For a given  $\mathbf{x} \in \mathcal{X}$ .

## 1) Softmax Response

**(SR)** [Hendrycks and Gimpel, 2017, Geifman and El-Yaniv, 2017]

ODIN with  $T = 1$  and  $\epsilon = 0$ .

## 2) Mahalanobis distance (MHLNB) [Lee et al., 2018]

Mahalanobis distance without input pre-processing and with the softmax output in place of the logits.

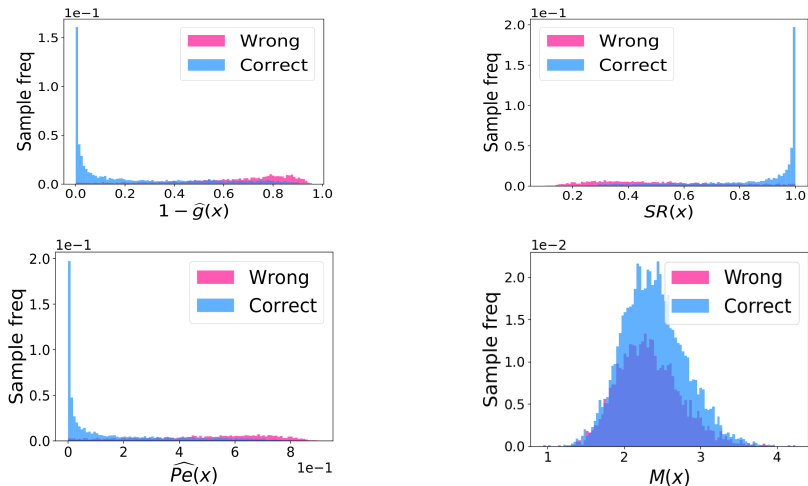
## TBB

- \* Temperature scaling,  $T = 1$
- \* Input pre-processing,  $\epsilon = 0$

## PBB

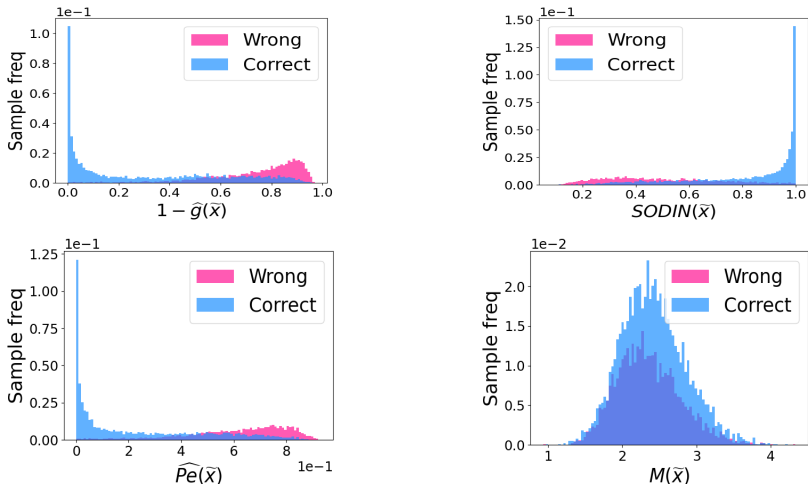
- \*  $D_\alpha, T_\alpha = 1$  and  $\epsilon_\alpha = 0.00035$
- \*  $D_\beta, T_\beta = 1.5$  and  $\epsilon_\beta = 0.00035$
- \* ODIN,  $T_{\text{ODIN}} = 1.3$  and  $\epsilon_{\text{ODIN}} = 0$
- \* MHLNB,  $T_{\text{MHLNB}} = 1$  and  $\epsilon_{\text{MHLNB}} = 0.0002$

# Discrimination performance for TBB

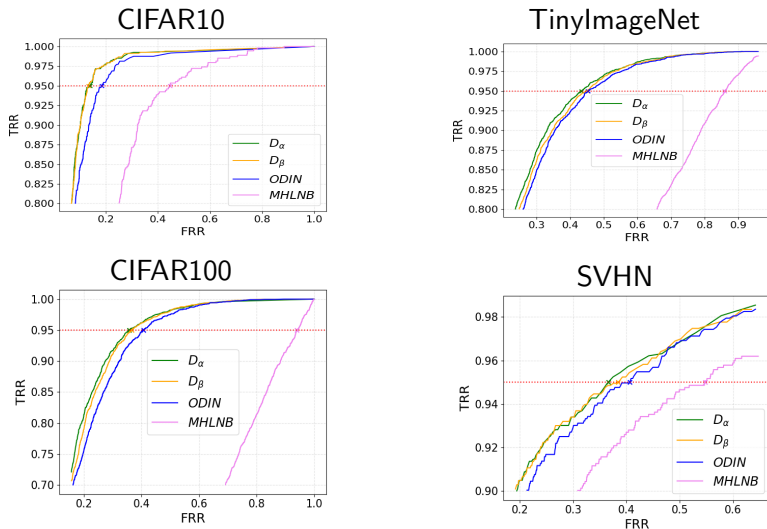


**Figure 1.** DOCTOR, SR and MHLNB to split data samples in TinyImageNet under TBB. Histograms for **wrongly classified samples** and **correctly classified samples**.

# Discrimination performance for PBB



**Figure 2.** DOCTOR, ODIN and MHLNB to split data samples in TinyImageNet under PBB. Histograms for **wrongly classified samples** and **correctly classified samples**.



**Figure 3.** ROC curves. Comparison between **DOCTOR**, **ODIN** and **MHLNB**. The red dashed line marks the 95% threshold of TRR.

# Overall Results: TBB & PBB

**Table 1.** Collection of the results in both **TBB** and **PBB**. For all methods, in TBB, we set  $T = 1$  and  $\epsilon = 0$ ; in PBB we set :  $\epsilon_\alpha = 0.00035$  and  $T_\alpha = 1$ ,  $\epsilon_\beta = 0.00035$  and  $T_\beta = 1.5$ ,  $\epsilon_{\text{ODIN}} = 0$  and  $T_{\text{ODIN}} = 1.3$ ,  $\epsilon_{\text{MHLNB}} = 0.0002$  and  $T_{\text{MHLNB}} = 1$ . In TBB for ODIN we report same results as in SR, since both methods coincide when  $T = 1$  and  $\epsilon = 0$ .

DATASET	METHOD	AUROC %		FRR % (95 % TRR)	
		TBB	PBB	TBB	PBB
CIFAR10 Acc. 95%	$D_\alpha$	<b>94</b>	<b>95.2</b>	<b>17.9</b>	13.9
	$D_\beta$	68.5	94.8	18.6	<b>13.4</b>
	ODIN	93.8	94.2	18.2	18.4
	SR	93.8	-	18.2	-
	MHLNB	92.2	84.4	30.8	44.6
CIFAR100 Acc. 78%	$D_\alpha$	<b>87</b>	<b>88.2</b>	40.6	<b>35.7</b>
	$D_\beta$	84.2	87.4	40.6	36.7
	ODIN	86.9	87.1	40.5	40.7
	SR	86.9	-	<b>40.5</b>	-
	MHLNB	82.6	50	66.7	94
TINY IMAGENET Acc. 63%	$D_\alpha$	<b>84.9</b>	<b>86.1</b>	<b>45.8</b>	<b>43.3</b>
	$D_\beta$	<b>84.9</b>	85.3	<b>45.8</b>	45.1
	ODIN	84.9	84.9	45.8	45.3
	SR	<b>84.9</b>	-	<b>45.8</b>	-
	MHLNB	78.4	59	82.3	86
SVHN Acc. 96%	$D_\alpha$	<b>92.3</b>	<b>93</b>	<b>38.6</b>	<b>36.6</b>
	$D_\beta$	92.2	92.8	39.7	38.4
	ODIN	92.3	92.3	38.6	40.7
	SR	<b>92.3</b>	-	<b>38.6</b>	-
	MHLNB	87.3	88	85.8	54.7
AMAZON FASHION Acc. 85%	$D_\alpha$	<b>89.7</b>	-	27.1	-
	$D_\beta$	<b>89.7</b>	-	<b>26.3</b>	-
	SR	87.4	-	50.1	-
AMAZON SOFTWARE Acc. 73%	$D_\alpha$	<b>68.8</b>	-	<b>73.2</b>	-
	$D_\beta$	<b>68.8</b>	-	<b>73.2</b>	-
	SR	67.3	-	86.6	-
IMDB Acc. 90%	$D_\alpha$	<b>84.4</b>	-	<b>54.2</b>	-
	$D_\beta$	<b>84.4</b>	-	54.4	-
	SR	83.7	-	61.7	-



# Misclassification Detection in Presence of OOD Samples

- ✦ DOCTOR is not tuned for OOD detection (differently from ODIN).
- ✦ We test ODIN and DOCTOR when one sample to reject out of five ( $\clubsuit$ ), three ( $\diamond$ ), or two ( $\spadesuit$ ) is OOD.

DATASET- In	DATASET- Out	AUROC %				FRR % (95 % TRR)			
		$D_\alpha$	$D_\beta$	ODIN	ENERGY	$D_\alpha$	$D_\beta$	ODIN	ENERGY
CIFAR10 $\clubsuit$	iSUN	<b>95.4</b> / 0.1	95.1 / 0.1	94.6 / 0.1	92.4 / 0	14 / 0.5	<b>13.5</b> / 0.4	17.2 / 0.3	32.2 / 0.1
	TINY (RES)	<b>95.2</b> / 0.1	94.9 / 0	94.6 / 0.1	92.3 / 0.1	<b>14</b> / 0.4	<b>14</b> / 0.5	17.8 / 0.4	32.2 / 0.1
CIFAR10 $\diamond$	iSUN	<b>95.5</b> / 0.1	95.3 / 0.1	94.9 / 0.1	92.9 / 0	14.4 / 0.6	<b>13.4</b> / 0.2	16.8 / 0.5	27 / 1
	TINY (RES)	<b>95.4</b> / 0.1	95 / 0.1	94.8 / 0.1	92.8 / 0	15 / 0.1	<b>14.8</b> / 0.7	17 / 0.5	28.8 / 1.9
CIFAR10 $\spadesuit$	iSUN	<b>95.6</b> / 0.1	<b>95.6</b> / 0	95.4 / 0	93.6 / 0.1	15.1 / 0.1	<b>13.6</b> / 0.5	16.1 / 0.2	25.1 / 0.2
	TINY (RES)	<b>95.5</b> / 0.1	95.2 / 0.1	95.1 / 0.1	93.5 / 0	<b>14.7</b> / 0.3	14.8 / 0.5	17.1 / 0.4	25.6 / 0.3

**Table 2.** Results in terms of *mean / standard deviation*.

# Takeaways from DOCTOR

- ✦ DOCTOR provides a flexible framework for misclassification error detection that applies to any pre-trained DNN classifier.
- ✦ We leverage information-theoretic tools to better discriminate between trusted and untrusted model predictions.
- ✦ Our method adapts to various scenarios depending on the level of information access of the DNN, uses only the pre-trained model.

## On-going work:

- ✦ Formalize statistical learning mechanisms that enable error detection and adaptation from few resources.
- ✦ Characterize their capabilities and limitations.
- ✦ Extension to semantic image segmentation, object detection and regression problems.

## Supplementary: Optimal (Oracle) Discriminator

- \*  $E \triangleq \mathbb{1}[Y \neq f_{\mathcal{D}_n}(\mathbf{X})]$  denotes the **error variable corresponding to  $f_{\mathcal{D}_n}$**
- \*  $\mathbf{x} \in \mathcal{X}$  and  $y \in \mathcal{Y}$  drawn from the unknown distribution  $p_{XY}$
- \*  $p_{XY}(\mathbf{x}, y) \equiv P_E(1)p_{XY|E}(\mathbf{x}, y|1) + P_E(0)p_{XY|E}(\mathbf{x}, y|0)$
- \*  $p_X(\mathbf{x}) \equiv P_E(1)p_{X|E}(\mathbf{x}|1) + P_E(0)p_{X|E}(\mathbf{x}|0)$
- \*  $\text{Pe}(\mathbf{x}) \triangleq \mathbb{E}[E(\mathbf{x})|\mathbf{x}] = 1 - P_{Y|X}(f_{\mathcal{D}_n}(\mathbf{x})|\mathbf{x})$  is the **probability of error classification w.r.t.  $P_{Y|X}$**

$$\begin{aligned} D(\mathbf{x}, \gamma) &= \mathbb{1}[p_{X|E}(\mathbf{x}|1) > \gamma \cdot p_{X|E}(\mathbf{x}|0)] \\ &= \mathbb{1}[P_{E|X}(1|\mathbf{x})P_E(0) > \gamma \cdot (1 - P_{E|X}(1|\mathbf{x}))P_E(1)] \\ &= \mathbb{1}[\text{Pe}(\mathbf{x})P_E(0) > \gamma \cdot (1 - \text{Pe}(\mathbf{x}))P_E(1)] \\ &= \mathbb{1}[\text{Pe}(\mathbf{x}) > \gamma' \cdot (1 - \text{Pe}(\mathbf{x}))], \end{aligned}$$

where  $\gamma' = \frac{P_E(1)}{P_E(0)}$ .



Geifman, Y. and El-Yaniv, R. (2017).

Selective classification for deep neural networks.

In Guyon, I., von Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4878–4887.



Hendrycks, D. and Gimpel, K. (2017).

A baseline for detecting misclassified and out-of-distribution examples in neural networks.

In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.



Lee, K., Lee, K., Lee, H., and Shin, J. (2018).

A simple unified framework for detecting out-of-distribution samples and adversarial attacks.

In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 7167–7177.



Liang, S., Li, Y., and Srikant, R. (2018).

Enhancing the reliability of out-of-distribution image detection in neural networks.

*In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.*

# IGOOD: An Information Geometry Approach to Out-of-Distribution Detection

Joint work with Eduardo D. C. Gomes, Florence Alberge  
and Pierre Duhamel



([https://openreview.net/pdf?id=mfwdY3U\\_9ea](https://openreview.net/pdf?id=mfwdY3U_9ea))

- We introduce IGEOOD, an effective method for detecting **Out-of-Distribution (OOD)** samples.
- IGEOOD applies to any pre-trained neural network, works under different degrees of access to the ML model, does not require OOD samples or assumptions on the OOD data but can also benefit (if available) from OOD samples.
- By building on the geodesic (**Fisher-Rao**) distance between the underlying data distributions, our discriminator combines confidence scores from the logits outputs and the learned features of a deep neural network.

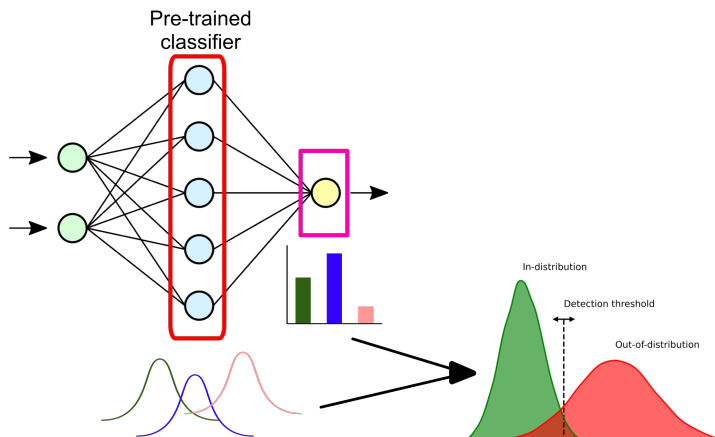
- Let  $\mathcal{X} \subseteq \mathbb{R}^d$  be the feature space and  $\mathcal{Y}$  a label space and let  $p_{\mathcal{X}\mathcal{Y}}$  be the underlying unknown probability density function (pdf) over  $\mathcal{X} \times \mathcal{Y}$ .
- In order to model the underlying problem, we introduce an artificial binary random variable  $Z \in \{0, 1\}$  indicating with  $z = 1$  that the test sample  $\mathbf{x}$  is OOD and  $z = 0$  otherwise.
- The open-world data can then be modeled as a *mixture* distribution  $p_{\mathcal{X}|Z}$  defined by

$$p_{\mathcal{X}|Z}(\mathbf{x}|z = 0) \triangleq p_{\mathcal{X}}(\mathbf{x}), \quad p_{\mathcal{X}|Z}(\mathbf{x}|z = 1) \triangleq q_{\mathcal{X}}(\mathbf{x}).$$

- The intrinsic difficulty arises from the fact that very little can be assumed about the unknown distributions  $p_{\mathcal{X}}$  and  $q_{\mathcal{X}}$ , in particular for out-of-distribution.
- **Alternative:** distance based criteria w.r.t an in-distribution probability reference.



# Statistical Model



**Figure:** We model the hidden layers' outputs as class conditional Gaussian distributions and the DNN's outputs as softmax probability distributions.

# Fisher-Rao Geodesic Distance

We propose an OOD detector based on the geodesic Fisher-Rao distance between probability density functions:

$$d_{\text{FR}}(q_{\theta}, q_{\theta}') \triangleq \inf_{\gamma} \int_0^1 \sqrt{\frac{d\gamma^{\top}(t)}{dt} \mathbf{G}(\gamma(t)) \frac{d\gamma(t)}{dt}} dt$$

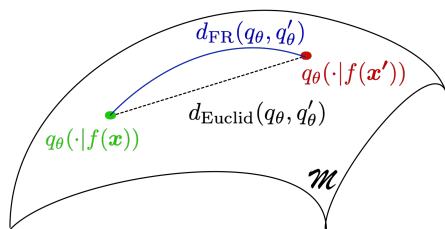


Figure: Illustration of the shortest path between distributions in a statistical manifold.

# IGOOD Score Using the Soft-Predictions

- **Igeood score using the softmax probability:** Let  $q_{\theta}(\cdot|f(\mathbf{x}))$  be the softmax probability distribution of the outputs. We can define the Fisher-Rao distance between softmax distributions as:

$$d_{\text{FR-Logits}}(q_{\theta}, q'_{\theta}) \triangleq 2 \arccos \left( \sum_{y \in \mathcal{Y}} \sqrt{q_{\theta}(y|f(\mathbf{x}))q_{\theta}(y|f(\mathbf{x}'))} \right)$$

- From which we derive our IGOOD score for the logits:

$$\text{FR}_0(\mathbf{x}) \triangleq \sum_{y \in \mathcal{Y}} d_{\text{FR-Logits}}(q_{\theta}(\cdot|f(\mathbf{x})), q_{\theta}(\cdot|\mu_y))$$

- Where  $\mu_y$  are the class conditional centroids given by:

$$\mu_y \triangleq \min_{\mu \in \mathbb{R}^{|\mathcal{Y}|}} \frac{1}{N_y} \sum_{\forall i: y_i=y} d_{\text{FR-Logits}}(q_{\theta}(\cdot|f(\mathbf{x}_i)), q_{\theta}(\cdot|\mu)).$$

# IGOOD Score Leveraging Latent Features

- **Igood score leveraging latent features:** For each layer, we model the features as a set of class-conditional Gaussian distributions with diagonal standard deviation matrix:

$$\mu_y^{(\ell)} = \frac{1}{N_y} \sum_{\forall i: y_i=y} f^{(\ell)}(\mathbf{x}_i)$$
$$\sigma^{(\ell)} = \text{diag} \left( \sqrt{\frac{1}{N} \sum_{y \in \mathcal{Y}} \sum_{\forall i: y_i=y} \left( f_j^{(\ell)}(\mathbf{x}_i) - \mu_{y,j}^{(\ell)} \right)^2} \right).$$

- We derive a confidence score by calculating the Fisher-Rao distance between the test sample  $\mathbf{x}$  and the closest class-conditional diagonal Gaussian distribution:

$$\text{FR}_\ell(\mathbf{x}) = \min_{y \in \mathcal{Y}} d_{\text{FR-Gauss}} \left( \left( \mathbf{x}, \sigma^{(\ell)} \right), \left( \mu_y^{(\ell)}, \sigma^{(\ell)} \right) \right).$$

- **Feature ensemble:** we combine the confidence scores of the logits and low-level features through a linear combination. If OOD data is available, we can also calculate  $\text{FR}'_{\ell}(\mathbf{x}; \boldsymbol{\mu}^{(\ell)'}, \boldsymbol{\sigma}^{(\ell)'})$  with OOD statistics, obtaining IGEOOD+:

$$\text{FR}(\mathbf{x}) \triangleq \alpha_0 \text{FR}_0(\mathbf{x}) + \sum_{\ell} \alpha_{\ell} \cdot \text{FR}_{\ell}(\mathbf{x}) + \alpha'_{\ell} \cdot \text{FR}'_{\ell}(\mathbf{x}).$$

- *Therefore, we have derived a unified OOD detection framework that combines a single distance for both the softmax outputs and the latent features of a neural network.*

# Experimental Setup

- The experimental setup follows the setting established by [1, 2, 4].
- We use two *pre-trained* deep neural networks architectures for image classification tasks: a Dense Convolutional Network (DenseNet-BC-100) and a Residual Neural Network (ResNet-34).
- *in-distribution data*: images from CIFAR-10, CIFAR-100 and SVHN datasets.
- *out-of-distribution data*: natural image examples from Tiny-ImageNet, LSUN, Describable Textures Dataset, Chars74K, Places365, iSUN and a synthetic dataset generated from Gaussian noise.

# Experimental Results

- The IGEOOD score increases the separation between in- and out-of-distribution data.

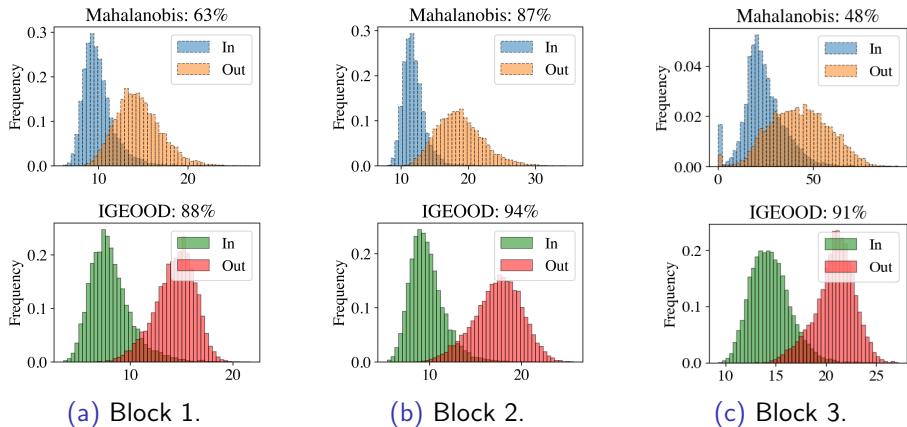


Figure: Histograms of the Mahalanobis and IGEOOD scores for the outputs of each hidden block of a DenseNet model.

# BLACK-BOX Results

**Table:** Average and standard deviation OOD detection performance across various OOD datasets for each model and in-distribution dataset in a BLACK-BOX setting. IGEOOD is compared to Baseline [1], ODIN [2], and Energy [3] methods.

Model	In-dist.	TNR at TPR-95%				AUROC			
		Baseline / ODIN / Energy / IGEOOD (ours)							
DenseNet	C-10	52.5 $\pm$ 16	<b>66.8</b> $\pm$ 20	65.3 $\pm$ 23	65.6 $\pm$ 23	91.8 $\pm$ 3.2	<b>92.8</b> $\pm$ 4.6	92.1 $\pm$ 5.3	92.3 $\pm$ 5.1
	C-100	15.9 $\pm$ 6.8	20.5 $\pm$ 9.5	20.3 $\pm$ 9.6	<b>20.7</b> $\pm$ 9.8	69.1 $\pm$ 15	71.6 $\pm$ 20	71.6 $\pm$ 20	<b>73.2</b> $\pm$ 17
	SVHN	68.4 $\pm$ 14	68.8 $\pm$ 20	70.2 $\pm$ 17	<b>72.1</b> $\pm$ 15	<b>92.3</b> $\pm$ 4.0	87.3 $\pm$ 14	90.1 $\pm$ 5.9	90.9 $\pm$ 5.3
ResNet	C-10	41.7 $\pm$ 16	51.9 $\pm$ 15	56.3 $\pm$ 13	<b>56.7</b> $\pm$ 13	89.6 $\pm$ 3.1	90.4 $\pm$ 3.1	90.4 $\pm$ 3.0	<b>90.5</b> $\pm$ 3.0
	C-100	15.0 $\pm$ 5.5	16.0 $\pm$ 6.3	16.3 $\pm$ 7.1	<b>16.4</b> $\pm$ 6.8	74.0 $\pm$ 1.9	75.2 $\pm$ 1.7	<b>75.5</b> $\pm$ 1.9	<b>75.5</b> $\pm$ 1.7
	SVHN	76.2 $\pm$ 7.8	77.7 $\pm$ 7.9	78.0 $\pm$ 7.9	<b>78.3</b> $\pm$ 8.0	<b>92.2</b> $\pm$ 2.9	91.4 $\pm$ 3.2	91.4 $\pm$ 3.2	91.7 $\pm$ 3.2
Average and Std.		44.9 $\pm$ 24	50.3 $\pm$ 24	51.1 $\pm$ 24	<b>51.6</b> $\pm$ 24	84.8 $\pm$ 9.5	84.8 $\pm$ 8.3	85.2 $\pm$ 8.4	<b>85.7</b> $\pm$ 8.0



# WHITE-BOX Results

- We increase the average TNR-95% by 11.8% and 2.5% with validation on OOD and adversarial data, respectively.

**Table:** Average and standard deviation of OOD detection performance for the WHITE-BOX settings. The abbreviation TNR-95%, C-10 and C-100 stands for TNR at TPR-95%, CIFAR-10 and CIFAR-100, respectively.

Model	In-dist.	Validation on OOD data		Validation on adversarial data	
		TNR-95%	AUROC	TNR-95%	AUROC
		Mahalanobis	IGEOOD+ (ours)	Mahalanobis	IGEOOD (ours)
DenseNet	C-10	76.6±31/ <b>92.6</b> ±14	92.1±12/ <b>98.4</b> ±3.0	75.9±30/ <b>77.9</b> ±29	91.7±12/ <b>94.0</b> ±9.0
	C-100	67.2±28/ <b>90.2</b> ±21	90.2±13/ <b>97.7</b> ±5.0	60.4±34/ <b>70.9</b> ±35	85.3±19/ <b>90.8</b> ±13
	SVHN	93.3±8.0/ <b>98.0</b> ±2.0	98.6±1.0/ <b>99.6</b> ±0.1	<b>93.7</b> ±10/92.2±9.0	<b>98.6</b> ±2.0/98.4±1.0
ResNet	C-10	82.5±23/ <b>91.6</b> ±16	96.5±4.0/ <b>98.4</b> ±3.0	<b>78.6</b> ±24/77.3±32	<b>95.3</b> ±6.0/90.0±15
	C-100	70.4±30/ <b>86.4</b> ±23	91.9±10/ <b>97.1</b> ±5.0	57.4±36/ <b>65.1</b> ±33	86.9±13/ <b>88.6</b> ±15
	SVHN	96.8±6.0/ <b>98.9</b> ±2.0	99.2±1.0/ <b>99.7</b> ±0.1	<b>96.3</b> ±8.0/93.6±14	<b>99.1</b> ±1.0/98.4±3.0
Average and Std.		81.1±11/ <b>92.9</b> ±4.0	94.8±4.0/ <b>98.5</b> ±1.0	77.0±15/ <b>79.5</b> ±10	92.8±5.4/ <b>93.4</b> ±3.9

# Further Comparison to the Literature

**Table:** TNR at TPR-95% (%) performance comparison in a WHITE-BOX setting considering the original results from [1,2,3,4]. Methods with an (\*) were tuned without OOD data.

	OOD dataset	CIFAR-10				CIFAR-100				SVHN			
		Mahalanobis [4] / Gram Matrix* [5] / DeConf-C* [6] / Res-Flow [7] / IGOOD / IGOOD+											
DenseNet	iSUN	95.3/99.0/	- / -	/97.7/ <b>99.8</b>	87.0/95.9/	- / -	/93.8/ <b>99.7</b>	<b>99.9</b> /99.4/	- / -	/98.3/ <b>99.9</b>			
	LSUN	97.2/99.5/99.4/98.2/98.5/ <b>99.9</b>			91.4/97.2/98.7/96.3/95.2/ <b>99.9</b>			<b>99.9</b> /99.5/	- /	<b>100</b> /97.1/ <b>99.9</b>			
	TinyImgNet	95.0/98.8/99.1/96.4/95.7/ <b>99.8</b>			86.6/95.7/98.6/93.0/94.5/ <b>99.5</b>			<b>99.9</b> /99.1/	- /	<b>100</b> /98.2/ <b>99.9</b>			
	SVHN/C-10	90.8/96.1/98.8/94.9/98.9/ <b>99.9</b>			82.5/89.3/95.9/84.9/93.3/ <b>99.6</b>			96.8/80.4/	- /	<b>99.0</b> /91.6/98.3			
ResNet	iSUN	97.8/99.3/	- / -	/97.2/ <b>99.9</b>	89.9/94.8/	- / -	/93.4/ <b>99.8</b>	99.7/99.4/	- / -	/99.8/ <b>100</b>			
	LSUN	98.8/99.6/	- /	99.0/98.4/ <b>100</b>	90.9/96.6/	- /	96.2/94.3/ <b>100</b>	<b>99.9</b> /99.6/	- /	<b>100</b> /99.7/ <b>99.9</b>			
	TinyImgNet	97.1/98.7/	- /	97.8/96.3/ <b>99.6</b>	90.9/94.8/	- /	94.6/90.1/ <b>99.6</b>	<b>99.9</b> /99.3/	- /	<b>100</b> /99.7/ <b>99.9</b>			
	SVHN/C-10	87.8/97.6/	- /	96.5/98.8/ <b>99.8</b>	91.9/80.8/	- /	93.0/91.6/ <b>99.7</b>	98.4/85.8/	- /	99.4/97.7/ <b>99.7</b>			

# Takeaways from IGEOOD

- IGEOOD provides a flexible framework for OOD detection that applies to any pre-trained DNN classifier.
- We leverage information geometry tools to better discriminate between probability distributions.
- Our method adapts to various scenarios depending on the level of information access of the DNN, uses only in-distribution samples but can also benefit (if available) of OOD samples.

## On-going work:

- Formalize hypothetical learning mechanisms that enable OOD generalization and adaptation.
- Characterize their capabilities and limitations.
- Extension to time-series and progressive distribution/model drifts.

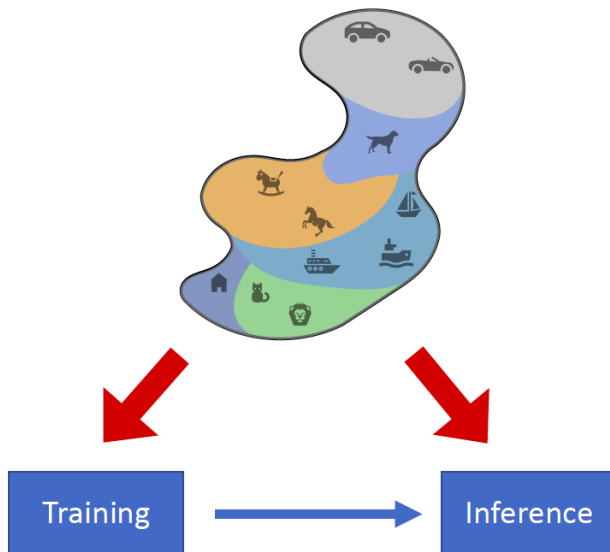
- [1] Dan Hendrycks & Kevin Gimpel. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. 2017.
- [2] Shiyu Liang et al. Enhancing The Reliability of Out-of-distribution Image Detection in Neural Networks. 2018.
- [3] Liu et al. Energy-based Out-of-distribution Detection. 2020.
- [4] Kimin L. et al. A simple unified framework for detecting out-of-distribution samples and adversarial attacks, 2018.
- [5] Sastry & Oore. Detecting out-of-distribution examples with Gram matrices, 2020.
- [6] Hsu et al. Generalized ODIN: Detecting out-of-distribution image without learning from out-of-distribution data, 2020.
- [7] Zisselman & Tamar. Deep residual flow for novelty detection, 2020.

# Adversarial Robustness via Fisher-Rao Regularization

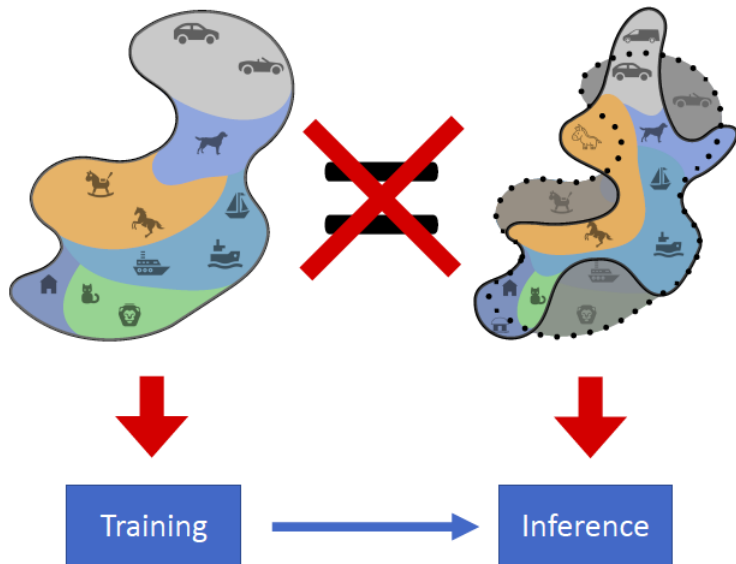
Joint work with Marine Picot, Francisco Messina, Malik Boudiaf,  
Fabrice Labeau, Ismail Ben Ayed

(<https://arxiv.org/abs/2106.06685>)

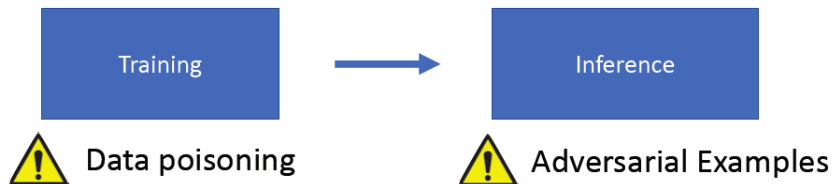
# Deep Neural Networks



# Deep Neural Networks

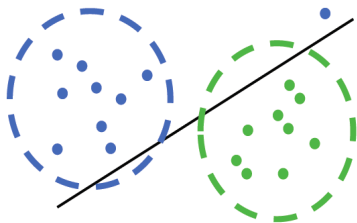
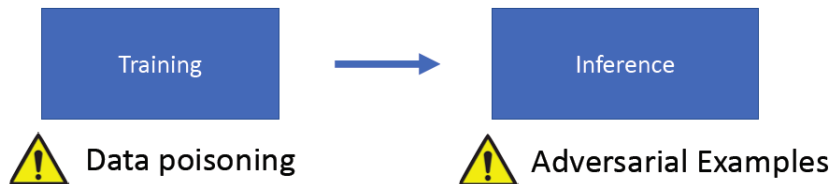


# Attacking Deep Neural Networks





# Attacking Deep Neural Networks

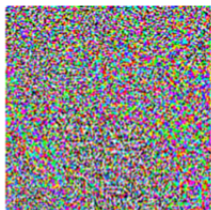


**Data poisoning:** modification of the boundaries

# Adversarial Examples



+ .007 ×



=



“panda”  
57.7% confidence

“nematode”  
8.2% confidence

“gibbon”  
99.3 % confidence

Figure: Building adversarial examples [Ian J Goodfellow et al. Arxiv 2014]

# Adversarial Examples

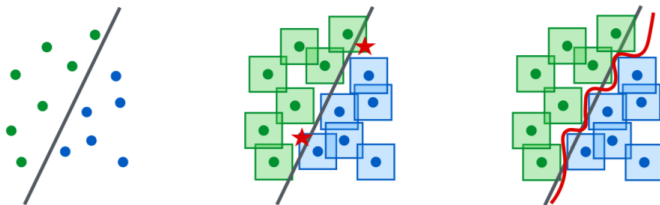


Figure: "Natural" vs "Adversarial" decision boundaries [A. Madry et al. ICLR 2018]

- Security



Glasses that fool face recognition [Mahmood Sharif et al. CCS 2016]



Let us consider the multi-class classification problem with:

- $\mathcal{X} \subseteq \mathbb{R}^n$  is the input space.
- $\mathcal{Y} = \{1, \dots, M\}$  is the label (concept) space.
- $q_\theta$  is the general classification model, parametrized by  $\theta \in \Theta$ .
- $P_e(\theta)$  is the error probability of the model parametrized by  $\theta \in \Theta$ .
- $\ell(\theta; \mathbf{x}, y)$  is the loss of the model parametrized by  $\theta$ , computed for the input  $(\mathbf{x}, y)$  and its expectation is the risk  $\mathcal{L}(\theta)$ .
- $\varepsilon$  is the maximal distortion allowed in the adversarial problem, according to a specific  $L^p$ -norm.
- $\mathbf{x}'$  refers to the adversarial version of any variable  $\mathbf{x}$ .

## Definition (Adversarial attacks)

The adversarial problem <sup>1</sup> is defined, according to  $L^p$ -norm, as:

$$\mathbf{x}^*(\mathbf{x}) \equiv \arg \min_{\mathbf{x}' \in [0,1]^n : \|\mathbf{x}' - \mathbf{x}\|_p < \epsilon} \|\mathbf{x}' - \mathbf{x}\|_p$$

s.t.  $f_\theta(\mathbf{x}') = t$

where

- $t$  is the target class or any class different from the original label  $y$ ,
- $\mathbf{x}' \in [0, 1]^n$  assures that  $\mathbf{x}^*(\mathbf{x})$  is close enough to the original image.

---

<sup>1</sup>Christian Szegedy et al. *Intriguing properties of neural networks* ICLR 2014.

# Fast Gradient Sign Method (FGSM)

Definition (FGSM Algorithm [Ian J Goodfellow et al. 2014])

$$\mathbf{x}' = \mathbf{x} + \alpha \operatorname{sgn}(\nabla_{\mathbf{x}} \ell(\theta; \mathbf{x}, y)),$$

where

- $(\mathbf{x}, y)$ : clean example
- $\mathbf{x}'$ : adversarial example
- $\operatorname{sgn}$  : the sign function
- $\nabla_{\mathbf{x}} \ell(\theta; \mathbf{x}, y)$ : the gradient w.r.t.  $\mathbf{x}$  of the loss function  $\ell(\theta; \mathbf{x}, y)$  evaluated at  $(\mathbf{x}, y)$
- $\alpha \leq \varepsilon$ : parameter controlling the magnitude of the perturbation.



## Definition (PGD Attack [A. Madry et al. ICLR 2018])

- It is the iterative extension of the FGSM method
- For a certain number of iterations  $k$ , we apply at each iteration  $i$ :

$$\mathbf{x}'^{(i+1)} = \mathbf{x}'^{(i)} + \delta \cdot \text{sgn} \left( \nabla_{\mathbf{x}} \ell(\theta; \mathbf{x}'^{(i)}, y) \right),$$

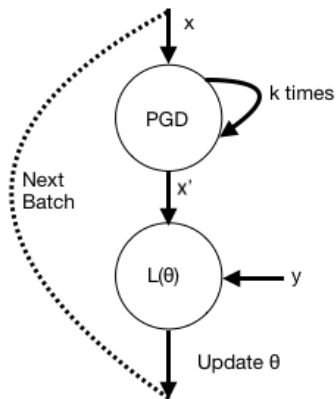
where  $\delta \leq \varepsilon$  is the noise norm at each step.

- $\mathbf{x}'^{(0)}$  is either equal to  $\mathbf{x}$  or  $\mathbf{x} + \boldsymbol{\eta}$  where  $\boldsymbol{\eta}$  is a random noise of maximum amplitude  $\varepsilon$ .
- To ensure that the  $L^p$ -norm constraint is met, at each iteration, we have to force:  $\|\mathbf{x}'^{(i)} - \mathbf{x}\|_p \leq \varepsilon$ .

# Adversarial Training

## Definition (Adversarial defense)

$$\theta^* \equiv \arg \min_{\theta} \mathbb{E}_{\mathbf{X}Y} \left[ \max_{\mathbf{x}': \|\mathbf{x}' - \mathbf{x}\|_{\infty} \leq \epsilon} \ell(\theta; \mathbf{x}', y) \right]$$



**Relaxation hypothesis:** We can approximate the max part with the generation of an adversarial example.

# Losses for Adversarial Training

Definition (Madry's method for defense [A. Madry et al. ICLR 2018])

Consider the adversarial cross-entropy loss (ACE):

$$\ell(\theta; \mathbf{x}', y) = -\log[q_{\theta}(y|\mathbf{x}')].$$

Definition (TRADES [Hongyang Zhang et al. ICML 2019])

Trade-off between natural and robust accuracies:

$$\ell(\theta; \mathbf{x}', y) = -\log[q_{\theta}(y|\mathbf{x})] + \lambda \cdot d_{KL}(q_{\theta}(y|\mathbf{x})\|q_{\theta}(y|\mathbf{x}')),$$

where  $\lambda$  is the hyperparameter controlling the trade-off between natural and adversarial accuracies.

**Robustness cannot be ensured against all adversarial (losses) attacks.** Can we derive an universal defense?

# Fisher-Rao Riemannian Geometry

## Definition (Fisher-Rao Distance (FRD))

- Given a family of probability distributions:

$$\mathcal{C} = \{q_\theta(\cdot|\mathbf{x}) : \mathbf{x} \in \mathcal{X}\}.$$

- Metric tensor (Fisher information):

$$G(\mathbf{x}) = \mathbb{E}_{Y \sim q_\theta(\cdot|\mathbf{x})} [\nabla_{\mathbf{x}} \log q_\theta(Y|\mathbf{x}) \nabla_{\mathbf{x}}^T \log q_\theta(Y|\mathbf{x})]$$

is positive definite for any  $\mathbf{x}$  and  $\theta \in \Theta$ .

- Infinitesimal squared length element:

$$ds^2 = \langle d\mathbf{x}, d\mathbf{x} \rangle_{G(\mathbf{x})} = d\mathbf{x}^T G(\mathbf{x}) d\mathbf{x}.$$

- The FRD between  $q_\theta(\cdot|\mathbf{x})$  and  $q_\theta(\cdot|\mathbf{x}')$  is:

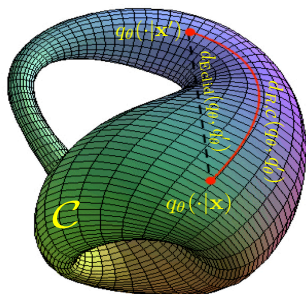
$$d_{R,\mathcal{C}}(q_\theta, q'_\theta) = \inf_{\gamma} \int_0^1 \sqrt{\frac{d\gamma^T(t)}{dt} G(\gamma(t)) \frac{d\gamma(t)}{dt}},$$

the inf is over all piecewise smooth curves.

- FRD is the length of the **geodesic** between  $(\mathbf{x}, \mathbf{x}')$  using  $G(\mathbf{x})$  as the metric tensor.



R. Rao and R. Fisher, 1956



## Definition (Fisher-Rao Distance (FRD))

We define the FIRE loss function as the trade-off between the natural cross-entropy and the expected Fisher-Rao distance between natural and adversarial probability distributions:

$$\ell_{\text{FIRE}}(\theta; \mathbf{x}, y) = -\log q_{\theta}(y|\mathbf{x}) + \lambda \cdot d_R^2(q_{\theta}(\cdot|\mathbf{x}), q_{\theta}(\cdot|\mathbf{x}')),$$

where  $\lambda$  is the hyperparameter controlling the trade-off between natural and adversarial performances with

$$d_R(q_{\theta}(\cdot|\mathbf{x}), q_{\theta}(\cdot|\mathbf{x}')) = 2 \arccos \left( \sum_{y \in \mathcal{Y}} \sqrt{q_{\theta}(y|\mathbf{x})q_{\theta}(y|\mathbf{x}')} \right).$$

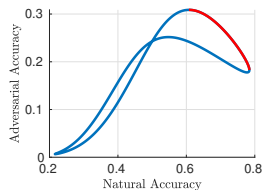
This metric has very interesting properties and **is related to well-known distances and Information divergences.**

# Comparison to the Kullback-Leiber Distance

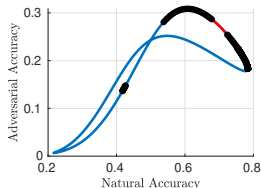
## Definition (Binary logistic regression)

Assume two equally likely classes  $\mathcal{Y} = \{-1, 1\}$  with conditional inputs given by  $\mathbf{x}|y \sim \mathcal{N}(y\mu, \Sigma)$ , and softmax probability

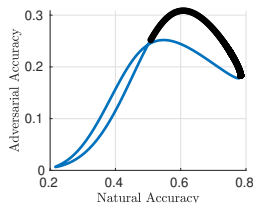
$$q_{\theta}(y|\mathbf{x}) = \frac{1}{1 + \exp(-y \theta^{\top} \mathbf{x})}.$$



(a) Pareto-optimal points



(b) TRADES



(c) FIRE

Figure: Plot of all possible pairs  $(1 - P_e(\theta), 1 - P'_e(\theta))$  for Gaussian model  $\varepsilon = 0.1$

# Experimental Set-Up

- **Datasets:** MNIST, CIFAR-10 - with and without additional data (AD) - , CIFAR-100
- **Model architecture:** CNNs, ResNet
- **Training procedure without AD:** Number of epochs : 100, batch size : 256, optimizer: SGD with a 0.9 momentum, and  $1.10^{-4}$  weight decay,  $l_r$ : 0.01 for MNIST, and to 0.1 for CIFAR-10 and CIFAR-100,  $l_r$  decay :divided by 10 at epochs 75 and 90.
- **Changes for AD simulations:** Number of epochs: 200,  $l_r$  decay: cosine.
- **Generation of adversarial examples:** PGD.
- **Additional data:** 500k additional images from 80M-TI<sup>1</sup>, selected such that the  $l_2$ -norm between those images and the images from CIFAR-10 are below a threshold.

---

<sup>1</sup>Images available at <https://github.com/yaircarmon/semisup-adv>

# Comparison between KL and Fisher-Rao Regularization

**Table:** Comparison between KL and Fisher-Rao based regularizer under white-box  $l_\infty$  threat model.

Defense	Dataset	$\epsilon$	Structure	Natural	AutoAttack	AA	RunTime
TRADES	MNIST	0.3	CNN	99.35	92.91	96.13	2h22
FIRE			CNN	99.13	94.06	96.59	2h06
TRADES	CIFAR-10	8/255	WRN-34-10	86.01	50.26	68.13	13h49
FIRE			WRN-34-10	85.42	52.22	68.82	11h00
TRADES	CIFAR-100	8/255	WRN-34-10	59.76	26.09	42.92	13h49
FIRE			WRN-34-10	60.71	27.63	44.17	11h10



# Comparison to SOTA Defense Mechanisms

**Table:** Test robustness on different datasets under white-box  $l_\infty$  attack. '\*' indicates models were retrained. '-' indicates the result is unavailable.

Defense	Dataset	$\epsilon$	Structure	Natural	AutoAttack	AA	Runtime
<b>Without Additional Data</b>							
Madry et al.	MNIST	0.3	CNN	98.53	88.50	93.51	2h03
Atzmon et al.			CNN	99.35	90.85	95.10	-
TRADES *			CNN	99.35	92.91	96.13	2h22
FIRE			CNN	99.14	94.06	96.60	2h06
Madry et al.	CIFAR-10	8/255	WRN-34-10	87.14	44.04	65.59	10h51
TRADES *			WRN-34-10	84.79	51.92	68.35	13h49
Self Adaptive			WRN-34-10	83.48	53.34	68.41	13h57
Overfitting *			WRN-34-10	86.85	51.74	69.29	42h01
FIRE	WRN-34-10	85.20	53.49	69.35	11h00		
Overfitting	CIFAR-100	8/255	RN-18	53.83	18.95	36.39	-
Overfitting*			WRN-34-10	59.01	27.07	43.04	42h08
FIRE			WRN-34-10	60.71	27.63	44.17	11h10
<b>With Additional Data Using 80M-TI</b>							
Pre-training	CIFAR-10	8/255	WRN-28-10	87.10	54.92	71.01	13h51
UAT			WRN-106-8	86.46	56.03	71.24	-
MART			WRN-28-10	87.50	56.29	71.89	10h22
RST-adv			WRN-28-10	89.70	59.53	74.61	22h12
FIRE			WRN-28-10	89.77	59.93	74.85	18h30

# Takeaways from FIRE

- FIRE is a novel method using tools from information geometry that encourages invariant softmax probabilities for natural and adversarial examples while maintaining high performances on natural samples.
- Theoretically, the optimization based on FIRE is well-behaved and gives all the desired Pareto-optimal points.
- Our empirical results showed that FIRE consistently enhances the robustness compared to TRADES.
- Compared to the state-of-the-art methods for adversarial defenses, FIRE increases the Average Accuracy (AA) while reducing the training time by 20%.

## On-going work:

- Our framework might be used to devise novel detection methods of adversarial examples.
- Characterize capabilities and limitations of potential attacks.
- Auditing mechanisms for ML models, based on partial statistical knowledge of the underlying distribution.

- 1 A Brief Overview of AI and Information Theory
  - The birth of AI and Deep Learning
  - Legacy of Shannon's work
  - Information, Uncertainty and Learning
- 2 Critical Problems in Safety AI
- 3 Overview of Recent Contributions to Safety AI
  - Detecting Misclassification Errors
  - Out-of-Distribution Detection
  - Adversarial Robustness
- 4 Discussion and Research Perspectives

## A long-lasting partnership:

- Learning is data compression.
- Concepts of data representations (e.g., encoders/decoders).
- Several information-based objectives (e.g., cross-entropy loss, mutual information,...), maximum information gain principle.
- Shannon entropy is a measure of randomness (or uncertainty).
- Minimum entropy principle is fundamental in statistical estimation and learning.

## Nonetheless, there is a long way to go:

- It is fundamentally important **to study other measures of information having more appropriate properties** from the viewpoint of its own learning problems.
- How we find an appropriate and universal way **to measure and to detect model uncertainty?**

From empirical evidence to information and knowledge:

- Researchers often have a tendency to fixate on **model performance metrics**, e.g., accuracy, but metrics **only tell part of the story** of a model's predictive decisions.
- It is important to understand what **drives a model to make predictions** (learning nature, not only imitation).

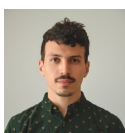
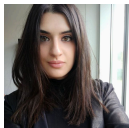
Uncertainty & robustness are critical problems in modern AI:

Models are often wrong, but **AI models that know when they are wrong are more useful**.

Better understanding of the information-theoretic link between:

- **Data** (source of empirical evidence),
- **Information** (the part that is unpredictable from the data),
- **Redundancy** (structure in data that provides the knowledge),
- **Knowledge** (the explanations of the complex world).

# Joint Work with PhD Students and Collaborators



- Federica Granese, Marine Picot, Eduardo D. C. Gomes, Francisco Messina, Malik Boudiaf, Marco Romanelli, Ismail Ben Ayed, Catuscia Palamidessi, Pierre Duhamel, Florence Alberge, Fabrice Labeau.

Thank you for your attention



**bpi**france

