Contrastive Learning for Novelty Detection

Jinwoo Shin

Korea Advanced Institute of Science and Technology (KAIST)

Joint work with Jihoon Tack* (KAIST), Sangwoo Mo* (KAIST), Jongheon Jeong (KAIST)

Paper draft is available at NeurIPS 2020 and https://arxiv.org/abs/2007.08176

Deep neural networks (DNNs) generalize well under the "seen" test distribution



For "unseen" distributions, DNNs often show astonishingly unexpected behaviors



Out-of-distribution

Unidentifiable



Adversarial samples

Is it possible to figure out whether a given sample is out-of-distribution (OOD)?



Practically, such an ability is indispensable for security-concerned systems



Autonomous drive



Secure authentication system

Is it possible to figure out whether a given sample is out-of-distribution (OOD)?

- **1.** How to learn a better representation $f(\cdot)$ more suitable for OOD detection?
- **2.** How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?



• : Out-of-distribution



Case 1: $f(\cdot) = a$ pre-trained classifier from a labeled dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

- **2.** How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?
- **Example 1** (the "Baseline" detector): Max-confidence based score [Hendrycks et al., 2017]



Case 1: $f(\cdot) = a$ pre-trained classifier from a labeled dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

- **2.** How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?
- Example 2: Mahalanobis-based confidence score [Lee et al., 2018]
 - Define a generative classifier $P(\mathbf{x}|y)$ from intermediate features

$$\mathbf{X} \Rightarrow \bigotimes_{i:y_i=c} \bullet \bullet \bullet \bigotimes_{i:y_i=c} f(\mathbf{x}_i), \quad \widehat{\mathbf{\Sigma}} = \frac{1}{N} \sum_{c} \sum_{i:y_i=c} (f(\mathbf{x}_i) - \widehat{\mu}_c) (f(\mathbf{x}_i) - \widehat{\mu}_c)^\top,$$

Lee, Lee, Lee & Shin. A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks, NeurIPS 2018.

Case 1: $f(\cdot) = a$ pre-trained classifier from a labeled dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

- **2.** How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?
- Example 2: Mahalanobis-based confidence score [Lee et al., 2018]
 - Define a generative classifier $P(\mathbf{x}|y)$ from intermediate features

$$s(\mathbf{x}) := \max_{c} \ - \left(f(\mathbf{x}) - \widehat{\mu}_{c}
ight)^{ op} \, \widehat{\mathbf{\Sigma}}^{-1} \left(f(\mathbf{x}) - \widehat{\mu}_{c}
ight)$$



 $P(\mathbf{x}|y)$: Generative, Mahalanobis-based



Lee, Lee, Lee & Shin. A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks, NeurIPS 2018.

Case 1: $f(\cdot) = a$ pre-trained classifier from a labeled dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

- **2.** How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?
- Example 2: Mahalanobis-based confidence score [Lee et al., 2018]
 (+) Near-perfect detection performances for "easy"-OODs, e.g., CIFAR-10 vs. LSUN
 (-) Still challenging on "harder"-OODs: e.g., CIFAR-10 vs. CIFAR-100 / One-class CIFAR-10

In-dist	000	TNR at TPR 95%	AUROC	Detection Acc.
(model)	000	Base	line / ODIN / Mahalanobis /	Ours
(iSUN	44.6/73.2/97.8/99.3	91.0 / 94.0 / 99.5 / 99.8	85.0 / 86.5 / 96.7 / 98.1
	LSUN (R)	49.8 / 82.1 / 98.8 / 99.6	91.0/94.1/99.7/99.9	85.3 / 86.7 / 97.7 / 98.6
CIEAD 10	LSUN (C)	48.6 / 62.0 / 81.3 / 89.8	91.9/91.2/96.7/97.8	86.3 / 82.4 / 90.5 / 92.6
(DecNet)	TinyImgNet (R)	41.0/67.9/97.1/98.7	91.0/94.0/99.5/99.7	85.1 / 86.5 / 96.3 / 97.8
(ResNet)	TinyImgNet (C)	46.4 / 68.7 / 92.0 / 96.7	91.4/93.1/98.6/99.2	85.4 / 85.2 / 93.9 / 96.1
	SVHN	50.5 / 70.3 / 87.8 / 97.6	89.9 / 96.7 / 99.1 / 99.5	85.1 / 91.1 / 95.8 / 96.7
	CIFAR-100	33.3 / 42.0 / 41.6 / 32.9	86.4 / 85.8 / 88.2 / 79.0	80.4 / 78.6 / 81.2 / 71.7

Results from [Sastry et al., 2020]

Lee, Lee, Lee & Shin. A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks, NeurIPS 2018. Sastry and Oore. Detecting Out-of-Distribution Examples with Gram Matrices, ICML 2020.



Case 1: $f(\cdot) = a$ pre-trained classifier from a labeled dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

2. How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?

Supervised representations can discriminate easy-OODs, but may not be enough for hard-OODs

-) Still challenging on narder -OODS: e.g., CIFAR-TU VS. CIFAR-TUU / One-class CIFAR-TU

In-dist	000	TNR at TPR 95%	AUROC	Detection Acc.
(model)	OOD TN iSUN 44.6 // LSUN (R) 49.8 // LSUN (C) 48.6 // TinyImgNet (R) 41.0 // TinyImgNet (C) 46.4 // SVHN 50.5 //	Basel	line / ODIN / Mahalanobis /	Ours
	iSUN	44.6 / 73.2 / 97.8 / 99.3	91.0 / 94.0 / 99.5 / 99.8	85.0 / 86.5 / 96.7 / 98.1
	LSUN (R)	49.8 / 82.1 / 98.8 / 99.6	91.0 / 94.1 / 99.7 / 99.9	85.3 / 86.7 / 97.7 / 98.6
CIEAD 10	LSUN (C)	48.6 / 62.0 / 81.3 / 89.8	91.9 / 91.2 / 96.7 / 97.8	86.3 / 82.4 / 90.5 / 92.6
(DecNet)	TinyImgNet (R)	41.0 / 67.9 / 97.1 / 98.7	91.0 / 94.0 / 99.5 / 99.7	85.1 / 86.5 / 96.3 / 97.8
(Residet)	TinyImgNet (C)	46.4 / 68.7 / 92.0 / 96.7	91.4 / 93.1 / 98.6 / 99.2	85.4 / 85.2 / 93.9 / 96.1
	SVHN	50.5 / 70.3 / 87.8 / 97.6	89.9 / 96.7 / 99.1 / 99.5	85.1 / 91.1 / 95.8 / 96.7
	CIFAR-100	33.3 / 42.0 / 41.6 / 32.9	86.4 / 85.8 / 88.2 / 79.0	80.4 / 78.6 / 81.2 / 71.7

Results from [Sastry et al., 2020]

Lee, Lee, Lee & Shin. A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks, NeurIPS 2018. Sastry and Oore. Detecting Out-of-Distribution Examples with Gram Matrices, ICML 2020.

Case 2: $f(\cdot) = a$ generative model from an unlabeled dataset $\mathcal{D} = \{\mathbf{x}_i\}$

- **1.** How to learn a better representation $f(\cdot)$ more suitable for OOD detection?
- **2.** How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?

Ideally, a good likelihood model $p(\mathbf{X})$ should also represent a good $s(\cdot)$



Case 2: $f(\cdot) = a$ generative model from an unlabeled dataset $\mathcal{D} = \{\mathbf{x}_i\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

2. How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?

Ideally, a good likelihood model $p(\mathbf{X})$ should also represent a good $s(\cdot)$

- Unfortunately, this is not the case at least for the current models
 - 1. They tend to be easily biased, e.g., to background statistics [Ren et al., 2019]
 - 2. Scaling up for a better likelihood model is usually much more challenging









(d) Train on ImageNet, Test on CIFAR-10 / CIFAR-100 / SVHN

Nalisnick et al. Do Deep Generative Models Know What They Don't Know. ICLR 2019. Ren et al. Likelihood Ratios for Out-of-Distribution Detection. NeurIPS 2019.

Case 2: $f(\cdot) = a$ generative model from an unlabeled dataset $\mathcal{D} = \{\mathbf{x}_i\}$

1. How to learn a better representation $f(\cdot)$ more suitable for OOD detection?

2. How to define a detection score $s(\cdot)$ that maximally utilizes $f(\cdot)$?











(d) Train on ImageNet, Test on CIFAR-10 / CIFAR-100 / SVHN

Nalisnick et al. Do Deep Generative Models Know What They Don't Know. ICLR 2019. Ren et al. Likelihood Ratios for Out-of-Distribution Detection. NeurIPS 2019.

- Learning unsupervised representation with self-supervision
- Example 1: Solving Jigsaw puzzles [Noroozi et al., 2016]



- Learning unsupervised representation with self-supervision
- Example 2: Predicting rotation angles (RotNet) [Gidaris et al., 2016]



Gidaris et al. Unsupervised Representation Learning by Predicting Image Rotations. ICLR 2018.

Hendrycks et al. (2019): RotNet improves novelty detection

$$\mathcal{L}_{SS}(x;\theta) = \frac{1}{4} \left[\sum_{r \in \{0^{\circ}, 90^{\circ}, 180^{\circ}, 270^{\circ}\}} \mathcal{L}_{CE}(\texttt{one_hot}(r), p_{\texttt{rot_head}}(r \mid R_r(x)); \theta) \right]$$

- $f(\cdot)$ is trained to predict the rotation angle {0°, 90°, 180°, 270°} applied to the input
- $s(\cdot)$ is defined to detect samples those failed to predict the applied rotations
- Intuition: Predicting rotations are harder to be transferred to OODs



Hendrycks, Mazeika, Kadavath and Song. Using Self-Supervised Learning Can Improve Model Robustness and Uncertainty. NeurIPS 2019.

Hendrycks et al. (2019): RotNet improves novelty detection

$$\mathcal{L}_{SS}(x;\theta) = \frac{1}{4} \left[\sum_{r \in \{0^{\circ}, 90^{\circ}, 180^{\circ}, 270^{\circ}\}} \mathcal{L}_{CE}(\texttt{one_hot}(r), p_{\texttt{rot_head}}(r \mid R_r(x)); \theta) \right]$$

w/ an external

• The proposed $s(\cdot)$ for RotNet largely advances AUROC on challenging one-class CIFAR-10

		-							OOD training data	
			OC-SVM	DeepSVDD	Geometric	RotNet	DIM	IIC	Supervised (OE)	Ours
	C	Airplane	65.6	61.7	76.2	71.9	72.6	68.4	87.6	77.5
		Automobile	40.9	65.9	84.8	94.5	52.3	89.4	93.9	96.9
		Bird	65.3	50.8	77.1	78.4	60.5	49.8	78.6	87.3
		Cat	50.1	59.1	73.2	70.0	53.9	65.3	79.9	80.9
CIFAR-10		Deer	75.2	60.9	82.8	77.2	66.7	60.5	81.7	92.7
classes		Dog	51.2	65.7	84.8	86.6	51.0	59.1	85.6	90.2
0140000		Frog	71.8	67.7	82.0	81.6	62.7	49.3	93.3	90.9
		Horse	51.2	67.3	88.7	93.7	59.2	74.8	87.9	96.5
		Ship	67.9	75.9	89.5	90.7	52.8	81.8	92.6	95.2
	C	Truck	48.5	73.1	83.4	88.8	47.6	75.7	92.1	93.3
		Mean	58.8	64.8	82.3	83.3	57.9	67.4	87.3	90.1

Hendrycks, Mazeika, Kadavath and Song. Using Self-Supervised Learning Can Improve Model Robustness and Uncertainty. NeurIPS 2019.

Hendrycks et al. (2019): RotNet improves novelty detection

$$\mathcal{L}_{SS}(x;\theta) = \frac{1}{4} \left[\sum_{r \in \{0^{\circ}, 90^{\circ}, 180^{\circ}, 270^{\circ}\}} \mathcal{L}_{CE}(\texttt{one_hot}(r), p_{\texttt{rot_head}}(r \mid R_r(x)); \theta) \right]$$

• T

Would this intuition still hold to more advanced SSL frameworks? ⇒ We examine a novelty detection from contrastive learning

	Airplane	65.6	61.7	76.2	71.9	72.6 68.4	87.6	77.5
	Automobile	40.9	65.9	84.8	94.5	52.3 89.4	93.9	96.9
	Bird	65.3	50.8	77.1	78.4	60.5 49.8	78.6	87.3
	Cat	50.1	59.1	73.2	70.0	53.9 65.3	79.9	80.9
CIFAR-10	Deer	75.2	60.9	82.8	77.2	66.7 60.5	81.7	92.7
classes	Dog	51.2	65.7	84.8	86.6	51.0 59.1	85.6	90.2
0100000	Frog	71.8	67.7	82.0	81.6	62.7 49.3	93.3	90.9
	Horse	51.2	67.3	88.7	93.7	59.2 74.8	87.9	96.5
	Ship	67.9	75.9	89.5	90.7	52.8 81.8	92.6	95.2
	Truck	48.5	73.1	83.4	88.8	47.6 75.7	92.1	93.3
	Mean	58.8	64.8	82.3	83.3	57.9 67.4	87.3	90.1

Hendrycks, Mazeika, Kadavath and Song. Using Self-Supervised Learning Can Improve Model Robustness and Uncertainty. NeurIPS 2019.

Learning representation that encodes the similarity between data points



Loss measured in the output space Examples: Colorization, Auto-Encoders Loss measured in the representation space Examples: TCN, CPC, Deep-InfoMax

Learning representation that encodes the similarity between data points



Examples: Colorization, Auto-Encoders

Loss measured in the representation space Examples: TCN, CPC, Deep-InfoMax

• A human prior defines positive & negative pairs w.r.t. a similarity score

$$\operatorname{score}(f(x), f(x^+)) \gg \operatorname{score}(f(x), f(x^-))$$

- Learning representation that encodes the similarity between data points
- We focus on SimCLR, one of representatives of modern CL [Chen et al., 2020]:
 - Attract (positive) the same instances with different data augmentations
 - Repel (negative) all the other different instances



Chen et al. A simple framework for contrastive learning of visual representations. ICML 2020.

- Learning representation that encodes the similarity between data points
- We focus on SimCLR, one of representatives of modern CL [Chen et al., 2020]:
 - Attract (positive) the same instances with different data augmentations



Chen et al. A simple framework for contrastive learning of visual representations. ICML 2020.

• SimCLR have largely closed the accuracy gap between un-/supervised learning



Summary: Contrasting Shifted Instances (CSI)

- We utilize the power of contrastive learning for OOD detection
- We further improve OOD detection by using shifted instances



Original



attract

• We train the representation via contrastive learning with shifted instances:



(SimCLR uses color jitter, random crop, horizontal flip, grayscale)

• We train the representation via contrastive learning with shifted instances:



26

- We train the representation via contrastive learning with shifted instances:
 - We found contrastive representation [Chen et al., 2020] is already good at OOD detection



: anchor

: attract

: repel

- We train the representation via **contrastive learning with shifted instances**:
 - We found **contrastive representation** [Chen et al., 2020] is already good at OOD detection
 - CSI further improves by **pushing the shifted samples** in addition to the different samples



- We train the representation via **contrastive learning with shifted instances**:
 - We found **contrastive representation** [Chen et al., 2020] is already good at OOD detection
 - CSI further improves by **pushing the shifted samples** in addition to the different samples
 - Additionally classify the shifting transformation



Contrasting Shifted Instances (CSI): Detection Score

• Detection score for **contrastively learned representation**:

• Further improving the detection score by **utilizing the shifting transformation**:

Contrasting Shifted Instances (CSI): Detection Score

- Detection score for contrastively learned representation:
 - The *cosine similarity* to the nearest training sample
 - The *norm* of the representation

• Further improving the detection score by **utilizing the shifting transformation**:

Contrasting Shifted Instances (CSI): Detection Score

- Detection score for contrastively learned representation:
 - The *cosine similarity* to the nearest training sample
 - The *norm* of the representation

- Further improving the detection score by utilizing the shifting transformation:
 - $s_{con-SI}(x, \{x_m\})$: ensemble the score $s_{con}(x; \{x_m\})$ over all shifting transformation
 - $s_{cls-SI}(x)$: confidence of the shifting transformation classifier

$$s_{\text{CSI}}(x; \{x_m\}) := s_{\text{con-SI}}(x; \{x_m\}) + s_{\text{cls-SI}}(x)$$

Contrasting Shifted Instances (CSI): OOD-ness

- **OOD-ness**: How to choose the shifting transformation?
 - The transformation that generates the most **OOD-like yet semantically meaningful samples**
 - We choose the transformation with the high OOD-ness (AUROC on vanilla SimCLR)



Contrasting Shifted Instances (CSI): Extension

- We also **extend CSI** for training <u>confidence-calibrated classifier</u> [Lee et al., 2018]:
 - Accurate on predicting label y when input x is in-distribution
 - Confidence $s_{sup}(x) \coloneqq \max_{y} p(y|x)$ of the classifier is well-calibrated

•: in-distribution *correct* sample •: in-distribution *in-correct* sample



 $S_{\sup}(\bigcirc) > S_{\sup}(\bigcirc)$

 $s_{\sup}(\bigcirc) > s_{\sup}(\bigodot)$

Contrasting Shifted Instances (CSI): Extension

- We also extend CSI for training confidence-calibrated classifier [Lee et al., 2018]:
 - Accurate on predicting label *y* when input *x* is in-distribution
 - Confidence $s_{sup}(x) \coloneqq \max_{y} p(y|x)$ of the classifier is well-calibrated



$$s_{\sup}(\bigcirc) > s_{\sup}(\bigcirc) \qquad s_{\sup}(\bigcirc) > s_{\sup}(\bigodot)$$

- We adapt the idea of CSI to the supervised contrastive learning (SupCLR):
 - SupCLR [Khosla et al., 2020] contrasts samples in *class-wise*, instead of in instance-wise
 - Similar to CSI, sup-CSI consider shifted instance as a different class's sample

Experiments: Unlabeled one-class OOD

CSI (ours)

• CSI achieves the state-of-the-art performance in all tested scenarios:

89.6

ResNet-18

• For unlabeled one-class OOD detection, outperforms prior methods under all classes

V	lethod	Network	Plane	Car	Bird	Cat	Deer	Dog	Frog	Horse	Ship	Truck	Mean
0	OC-SVM* [64]	_	65.6	40.9	65.3	50.1	75.2	51.2	71.8	51.2	67.9	48.5	58.8
D	DeepSVDD* [60]	LeNet	61.7	65.9	50.8	59.1	60.9	65.7	67.7	67.3	75.9	73.1	64.8
Ą	noGAN* [63]	DCGAN	67.1	54.7	52.9	54.5	65.1	60.3	58.5	62.5	75.8	66.5	61.8
0	CGAN* [55]	OCGAN	75.7	53.1	64.0	62.0	72.3	62.0	72.3	57.5	82.0	55.4	65.7
G	eom* [17]	WRN-16-8	74.7	95.7	78.1	72.4	87.8	87.8	83.4	95.5	93.3	91.3	86.0
R	lot* [27]	WRN-16-4	71.9	94.5	78.4	70.0	77.2	86.6	81.6	93.7	90.7	88.8	83.3
R	ot+Trans* [27]	WRN-16-4	77.5	96.9	87.3	80.9	92.7	90.2	90.9	96.5	95.2	93.3	90.1
G	60AD* [<mark>2</mark>]	WRN-10-4	77.2	96.7	83.3	77.7	87.8	87.8	90.0	96.1	93.8	92.0	88.2
R	lot [27]	ResNet-18	78.3 ± 0.2	$94.3{\scriptstyle \pm 0.3}$	$86.2{\scriptstyle \pm 0.4}$	$80.8{\scriptstyle \pm 0.6}$	$89.4{\scriptstyle \pm 0.5}$	$89.0{\scriptstyle \pm 0.4}$	$88.9{\scriptstyle \pm 0.4}$	$95.1{\scriptstyle \pm 0.2}$	$92.3{\scriptstyle \pm 0.3}$	$89.7{\scriptstyle\pm0.3}$	88.4
R	ot+Trans [27]	ResNet-18	80.4 ± 0.3	$96.4{\scriptstyle \pm 0.2}$	$85.9{\scriptstyle \pm 0.3}$	$81.1{\scriptstyle \pm 0.5}$	$91.3{\scriptstyle \pm 0.3}$	$89.6{\scriptstyle \pm 0.3}$	$89.9{\scriptstyle \pm 0.3}$	$95.9{\scriptstyle \pm 0.1}$	$95.0{\scriptstyle \pm 0.1}$	$92.6{\scriptstyle\pm0.2}$	89.8
G	GOAD [2]	ResNet-18	75.5 ± 0.3	$94.1{\scriptstyle\pm0.3}$	$81.8{\scriptstyle\pm0.5}$	72.0 ± 0.3	$83.7{\pm}0.9$	84.4 ± 0.3	$82.9{\scriptstyle\pm0.8}$	$93.9{\scriptstyle\pm0.3}$	$92.9{\scriptstyle\pm0.3}$	89.5±0.2	85.1
С	CSI (ours)	ResNet-18	89.9 ±0.1	$99.1{\scriptstyle\pm0.0}$	$93.1{\scriptstyle\pm0.2}$	$\pmb{86.4}{\scriptstyle\pm0.2}$	$\textbf{93.9}{\scriptstyle \pm 0.1}$	93.2±0.2	$95.1{\scriptstyle \pm 0.1}$	$98.7{\scriptstyle\pm0.0}$	$97.9{\scriptstyle \pm 0.0}$	95.5±0.1	94.3
	(b) One-cla	ss CIFAR-	100 (sup	er-class	5)			(c) One	-class In	nageNe	t-30		
	Method	Netv	vork	AURO	DC	Method	ł			Net	work	AUR	C
	OC-SVM* [54] -		63.1		Rot* [2	27]			Res	Net-18	65.3	3
	Geom* [17]	WRI	N-16-8	78.7	7	Rot+Tr	ans^* [2]	7]		Res	Net-18	77.9	9
	Rot [27]	ResN	Net-18	77.7	7	Rot+A	ttn* [27]		Res	Net-18	81.0	6
	Rot+Trans [2	[7] ResN	Net-18	79.8	3	Rot+Tr	ans+At	tn* [27]		Res	Net-18	84.8	8
	GOAD [2]	ResN	Net-18	74.5	5	Rot+Tr	ans+At	tn+Resi	ze* [27	Res	Net-18	85.7	7

CSI (ours)

ResNet-18

91.6

(a) One-class CIFAR-10

Experiments: Unlabeled multi-class OOD

CSI (ours)

ResNet-18

90.5±0.1

- CSI achieves the state-of-the-art performance in all tested scenarios:
 - For unlabeled multi-class OOD detection, outperforms prior methods under all OOD datasets

				$CIFAR10 \rightarrow$							
Method		Netwo	ork	SVHN	LSUN	ImageNet	LSU	JN (FIX)	ImageNet (FIX)	CIFAR-100	Interp.
Likelihood*		Pixel	CNN++	8.3	-	64.2	-		-	52.6	52.6
Likelihood*		Glow		8.3	-	66.3	-		-	58.2	58.2
Likelihood*		EBM		63.0	-	-	-		-	-	70.0
Likelihood Ratio*	[55]	Pixel	CNN++	91.2	-	-	-		-	-	-
Input Complexity*	* [61]	Pixel	CNN++	92.9	-	58.9	-		-	53.5	-
Input Complexity*	* [61]	Glow		95.0	-	71.6	-		-	73.6	-
Rot [25]		ResN	et-18	97.6±0.2	89.2±0.7	90.5±0.3	77.7	7±0.3	83.2±0.1	79.0 ± 0.1	64.0±0.3
Rot+Trans [25]		ResN	et-18	$97.8{\scriptstyle\pm0.2}$	$92.8{\scriptstyle \pm 0.9}$	$94.2{\pm}0.7$	81.6	5 ± 0.4	$86.7{\scriptstyle\pm0.1}$	82.3 ± 0.2	$68.1{\scriptstyle\pm0.8}$
GOAD [2]		ResN	et-18	$96.3{\scriptstyle \pm 0.2}$	$89.3{\scriptstyle\pm1.5}$	$91.8{\scriptstyle\pm1.2}$	78.8	3 ± 0.3	$83.3{\pm}0.1$	77.2 ± 0.3	59.4 ± 1.1
CSI (ours)		ResN	et-18	$99.8{\scriptstyle\pm0.0}$	$97.5{\scriptstyle\pm0.3}$	$97.6{\scriptstyle\pm0.3}$	90.3	3 ±0.3	93.3 ± 0.1	$89.2{\scriptstyle\pm0.1}$	$\textbf{79.3}{\scriptstyle \pm 0.2}$
				(b)) Unlabe	led Image	eNet-	-30			
							Imag	eNet-30 -	\rightarrow		
Method	Netwo	ork	CUB-2	00 Dog	gs Pe	ets Flow	vers	Food-10	1 Places-365	Caltech-256	DTD
Rot [25]	Rot [25] ResNet-18 76.5		76.5 ± 0	.7 77.2	±0.5 70.0	0±0.5 87.2	2 ± 0.2	72.7 ± 1.5	5 52.6 ±1.4	$70.9{\scriptstyle\pm0.1}$	$89.9{\scriptstyle\pm0.5}$
Rot+Trans [25]	ResNe	et-18	74.5 ± 0	.5 77.8 =	±1.1 70.0	0±0.8 86.3	± 0.3	71.6 ± 1.4	53.1±1.7	$70.0{\pm}0.2$	$89.4{\scriptstyle \pm 0.6}$
GOAD [2] ResNet-18		et-18	71.5 ± 1	.4 74.3	±1.6 65.5	5±1.3 82.8	± 1.4	68.7 ± 0.7	51.0±1.1	67.4 ± 0.8	$87.5{\pm}0.8$

97.1±0.1 **85.2**±0.2 **94.7**±0.4 **89.2**±0.3

78.3±0.3

87.1±0.1

96.9±0.1

(a) Unlabeled CIFAR-10

Experiments: Labeled multi-class OOD

- CSI achieves the state-of-the-art performance in all tested scenarios:
 - For labeled multi-class OOD detection, outperforms prior methods under all OOD datasets

				$CIFAR10 \rightarrow$								
Train method	Test acc.	ECE	SVHN	LSUN	ImageNet	LSUN (FIX)	ImageNet (FIX)	CIFAR100	Interp.			
Cross Entropy	$93.0{\pm}0.2$	$6.44{\pm}0.2$	88.6 ± 0.9	$90.7{\pm}0.5$	$88.3{\pm}0.6$	87.5 ± 0.3	87.4 ± 0.3	85.8 ± 0.3	$75.4{\pm}0.7$			
SupCLR [30]	$93.8{\scriptstyle\pm0.1}$	5.56 ± 0.1	$97.3{\scriptstyle \pm 0.1}$	$92.8{\scriptstyle \pm 0.5}$	$91.4{\scriptstyle\pm1.2}$	91.6±1.5	$90.5{\scriptstyle\pm0.5}$	88.6 ± 0.2	$75.7{\scriptstyle\pm0.1}$			
CSI (ours)	$94.8{\scriptstyle\pm0.1}$	4.40 ± 0.1	$96.5{\scriptstyle\pm0.2}$	$96.3{\scriptstyle \pm 0.5}$	96.2 ± 0.4	92.1 ± 0.5	$92.4{\scriptstyle\pm0.0}$	90.5 ± 0.1	$78.5{\scriptstyle \pm 0.2}$			
CSI-ens (ours)	$96.1{\scriptstyle \pm 0.1}$	$\textbf{3.50}{\scriptstyle \pm 0.1}$	$97.9{\scriptstyle \pm 0.1}$	$97.7{\scriptstyle\pm0.4}$	97.6±0.3	93.5 ± 0.4	$94.0{\scriptstyle\pm0.1}$	92.2 ± 0.1	$80.1{\pm}0.3$			

(a) Labeled CIFAR-10

(b) Labeled ImageNet-30

				ImageNet-30 \rightarrow							
Train method	Test acc.	ECE	CUB-200	Dogs	Pets	Flowers	Food-101	Places-365	Caltech-256	DTD	
Cross Entropy	94.3	5.08	88.0	96.7	95.0	89.7	79.8	90.5	90.6	90.1	
SupCLR [30]	96.9	3.12	86.3	95.6	94.2	92.2	81.2	89.7	90.2	92.1	
CSI (ours)	97.0	2.61	93.4	97.7	96.9	96.0	87.0	92.5	91.9	93.7	
CSI-ens (ours)	97.8	2.19	94.6	98.3	97.4	96.2	88.9	94.0	93.2	97.4	

Experiments: Ablation study

- We verified the effectiveness of **shifting transformation selection scheme**
 - Higher OOD-ness valued transformation leads to higher detection performance

					@ 1				
(a) Original	(b) Cut	out	(c) Sobel	(d) N	loise	(e) Blu	ır	(f) Perm	(g) Rotate
			Cutout	Sobel	Noise	Blur	Perm	Rotate	
	OOD-	ness	79.5	69.2	74.4	76.0	83.8	85.2	
				Ţ	I I High∉ ♦	er OOD	-ness -	→ Higher	performance
	Base		Cutout	Sobel	Noise	Blur	Perm	Rotate	
	87.9	+Align +Shift	a 84.3 88.5	85.0 88.3	85.5 89.3	88.0 89.2	73.1 90.7	76.5 94.3	

Experiments: Ablation study

- We verified the effectiveness of **shifting transformation selection scheme**
 - Higher OOD-ness valued transformation leads to higher detection performance
 - Our method works on rotation-invariant datasets *i.e., rotation is not shifting transformation*



(a) O	OD-ness		(b) AUROC				
Rot.	Noise	Base	Base CSI(R)				
50.6	75.7	70.3	65.9	80.1			

Experiments: Ablation study

- We verified the effectiveness of **shifting transformation selection scheme** •
 - Higher OOD-ness valued transformation leads to higher detection performance
 - Our method works on rotation-invariant datasets *i.e., rotation is not shifting transformation*

(a) OC	DD-ness		(b) AURC)C
Rot.	Noise	Base	CSI(R)	CSI(N)
50.6	75.7	70.3	65.9	80.1

Each of the proposed components is complementary for AUROC

	(a) Trainin	ig object	tive		(b) Detection score					
	SimCLR	Con.	Cls.	AUROC		Con.	Cls.	Ensem.	AUROC	
$\mathcal{L}_{\texttt{SimCLR}}$ (2)	\checkmark	_	-	87.9	s_{con} (6)	\checkmark	-	_*	91.3	
$\mathcal{L}_{\texttt{con-SI}}$ (3)	\checkmark	\checkmark	-	91.6	$s_{\texttt{con-SI}}(7)$	\checkmark	-	\checkmark	93.3	
$\mathcal{L}_{\texttt{cls-SI}}$ (4)	-	-	\checkmark	88.6	$s_{\texttt{cls-SI}}$ (8)	-	\checkmark	\checkmark	93.8	
$\mathcal{L}_{\texttt{CSI}}$ (5)	\checkmark	\checkmark	\checkmark	94.3	s_{CSI} (9)	\checkmark	\checkmark	\checkmark	94.3	

Conclusion

- We propose Contrasting Shifted Instances (CSI) for OOD detection
 - We extend the power of contrastive learning for OOD detection
 - We further improve the OOD detection by utilizing shifting transformations

CSI shows outstanding performance under various OOD detection scenarios

 We believe CSI would guide various future directions in OOD detection & selfsupervised learning as an important baseline