
Deep Mahalanobis Detector

2020. 02. 06

Soonwoo Kwon

soonwoo.kwon@kaist.ac.kr

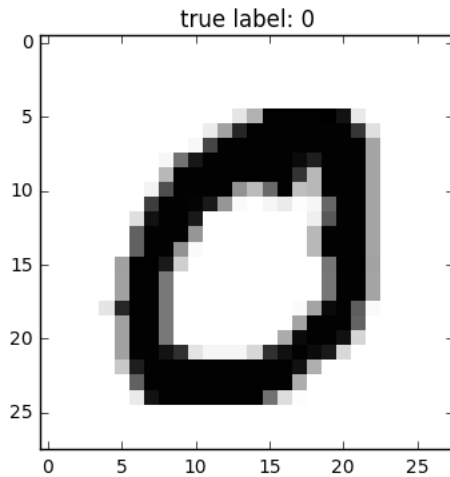
Table of Contents

- Concept of OOD (Out-Of-Distribution) detection
- **Deep Mahalanobis Detector**: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks , NIPS 2018

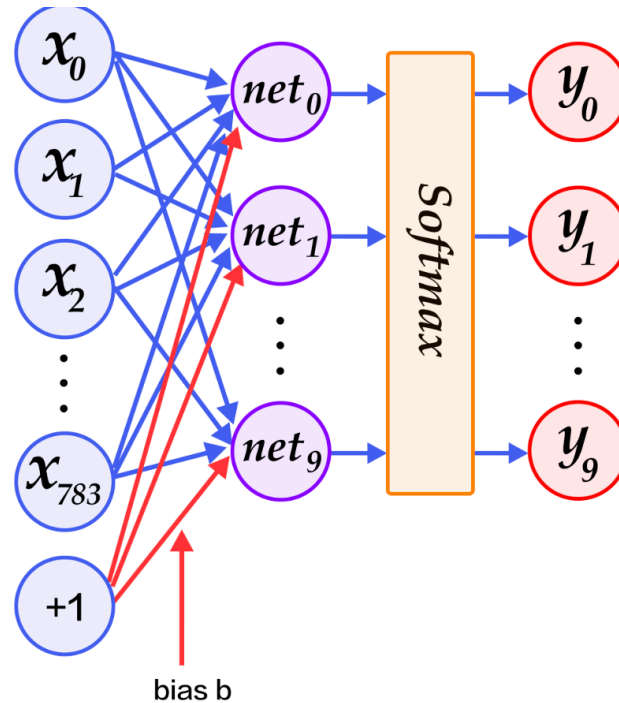
Table of Contents

- Concept of OOD (Out-Of-Distribution) detection
- Deep Mahalanobis Detector: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks , NIPS 2018

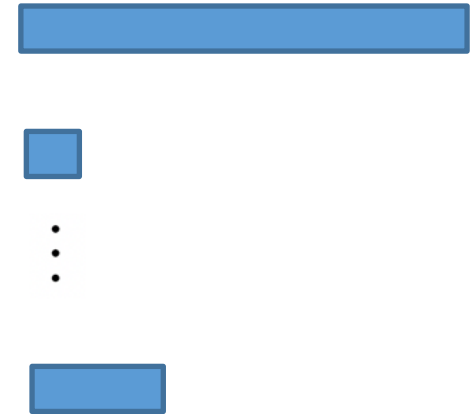
Motivation



Test Image
(MNIST)

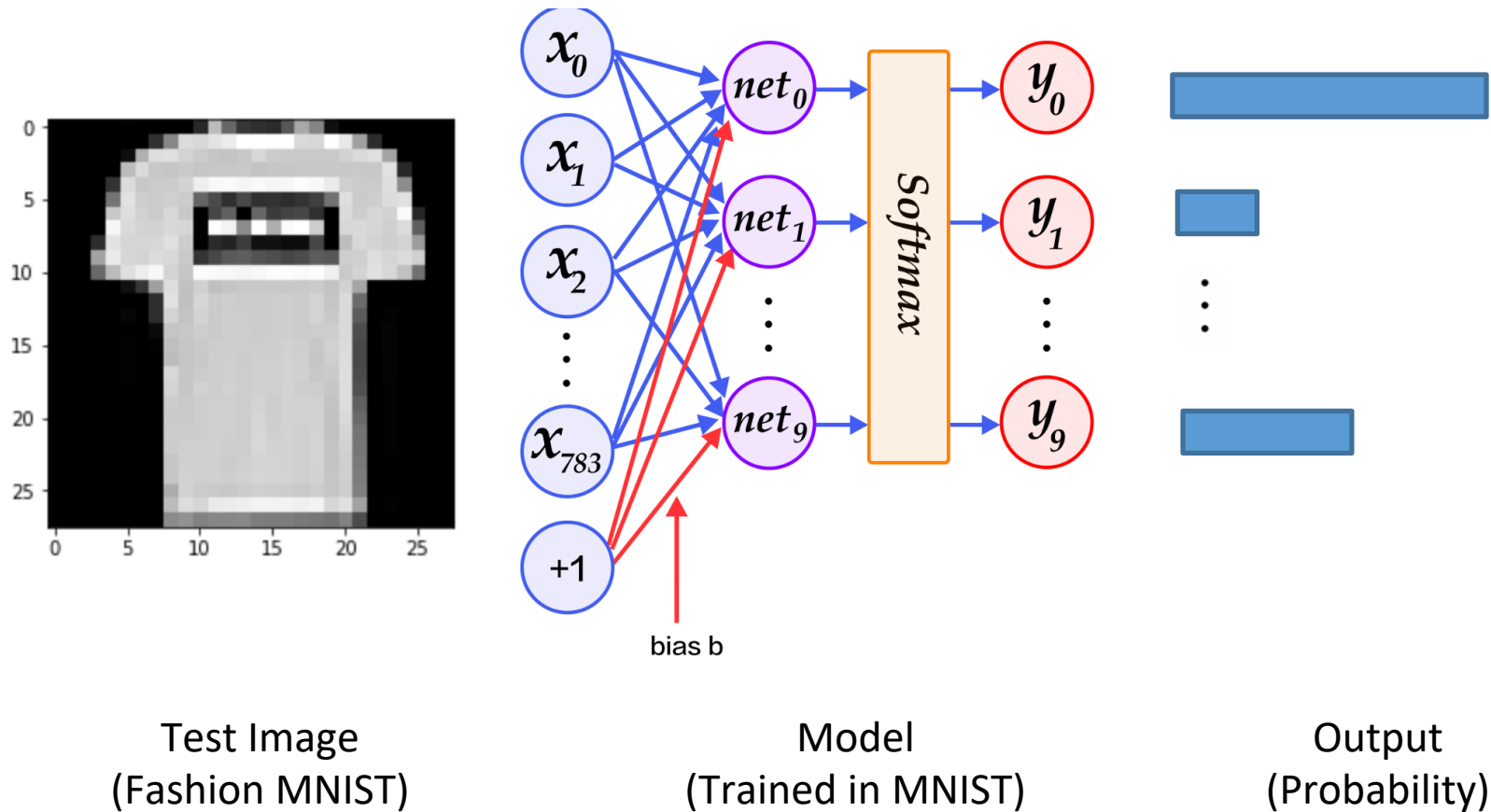


Model
(Trained on MNIST)



Output
(Probability)

Motivation



Still our model makes confident prediction..!

OOD(Out-Of-Distribution)?

- **In-distribution:** distribution of training samples
- **(OOD) Out-of-distribution**

Table of Contents

- Concept of OOD detection
- **Deep Mahalanobis Detector**: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks , NIPS 2018

Deep Mahalanobis Detector

- can use pre-trained model without any modification in test time
- is not only **simple** but also **shows good performance on OOD detection**
- can **detect** both **OOD samples** and **adversarial attack samples**
- can be extended to be utilized in class-incremental learning

Method

- Pretrained Softmax classifier:

$$P(y = c|\mathbf{x}) = \frac{\exp(\mathbf{w}_c^\top f(\mathbf{x}) + b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top f(\mathbf{x}) + b_{c'})}$$

(where, \mathbf{x} : input, c : class, $f(\cdot)$: penultimate layer of DNN)

- To get generative classifier,
 - **Suppose that class conditional distribution follows the multivariate Gaussian distribution.**

$$P(f(\mathbf{x})|y = c) = \mathcal{N}(f(\mathbf{x})|\mu_c, \Sigma)$$

empirical mean

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} f(\mathbf{x}_i),$$

empirical (tied) covariance

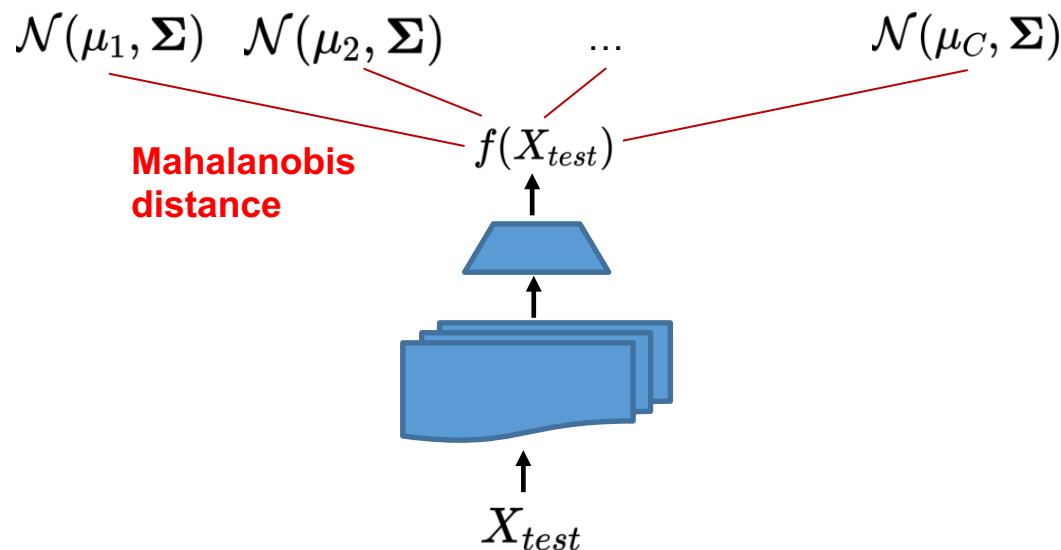
$$\hat{\Sigma} = \frac{1}{N} \sum_c \sum_{i:y_i=c} (f(\mathbf{x}_i) - \hat{\mu}_c)(f(\mathbf{x}_i) - \hat{\mu}_c)^\top$$

Method

- **Class conditional distribution** follows the **multivariate Gaussian distribution**.

$$P(f(\mathbf{x})|y = c) = \mathcal{N}(f(\mathbf{x})|\mu_c, \Sigma)$$

- Then, calculate **Mahalanobis distance** from $f(X_{test})$ to $\mathcal{N}(\mu_c, \Sigma)$, $c \in \{1, \dots, C\}$



- **Confidence score:** $M(\mathbf{x}) = \max_c - (f(\mathbf{x}) - \hat{\mu}_c)^\top \hat{\Sigma}^{-1} (f(\mathbf{x}) - \hat{\mu}_c)$
- If **confidence score** of X_{test} < **threshold**: X_{test} is predicted as **OOD**
- Classification: $\hat{y}(\mathbf{x}) = \arg \min_c (f(\mathbf{x}) - \hat{\mu}_c)^\top \hat{\Sigma}^{-1} (f(\mathbf{x}) - \hat{\mu}_c)$

Method

- Input pre-processing

$$\hat{\mathbf{x}} = \mathbf{x} + \varepsilon \text{sign}(\nabla_{\mathbf{x}} M(\mathbf{x})) = \mathbf{x} - \varepsilon \text{sign}\left(\nabla_{\mathbf{x}}(f(\mathbf{x}) - \hat{\mu}_{\hat{c}})^{\top} \hat{\Sigma}^{-1} (f(\mathbf{x}) - \hat{\mu}_{\hat{c}})\right)$$

Noise is generated to **increase** the proposed **confidence score**

- Feature ensemble
 - Measuring and combining the confidence scores from not only the final features, but also the other low-level features in DNNs

Algorithm 1 Computing the Mahalanobis distance-based confidence score.

Input: Test sample \mathbf{x} , weights of logistic regression detector α_ℓ , noise ε and parameters of Gaussian distributions $\{\hat{\mu}_{\ell,c}, \hat{\Sigma}_\ell : \forall \ell, c\}$

Initialize score vectors: $\mathbf{M}(\mathbf{x}) = [M_\ell : \forall \ell]$

for each layer $\ell \in 1, \dots, L$ **do**

Find the closest class: $\hat{c} = \arg \min_c (f_\ell(\mathbf{x}) - \hat{\mu}_{\ell,c})^\top \hat{\Sigma}_\ell^{-1} (f_\ell(\mathbf{x}) - \hat{\mu}_{\ell,c})$

Add small noise to test sample: $\hat{\mathbf{x}} = \mathbf{x} - \varepsilon \text{sign} \left(\nabla_{\mathbf{x}} (f_\ell(\mathbf{x}) - \hat{\mu}_{\ell,\hat{c}})^\top \hat{\Sigma}_\ell^{-1} (f_\ell(\mathbf{x}) - \hat{\mu}_{\ell,\hat{c}}) \right)$

Computing confidence score: $M_\ell = \max_c - (f_\ell(\hat{\mathbf{x}}) - \hat{\mu}_{\ell,c})^\top \hat{\Sigma}_\ell^{-1} (f_\ell(\hat{\mathbf{x}}) - \hat{\mu}_{\ell,c})$

end for

return Confidence score for test sample $\sum_\ell \alpha_\ell M_\ell$

Experimental Result

Method	Feature ensemble	Input pre-processing	TNR at TPR 95%	AUROC	Detection accuracy	AUPR in	AUPR out
Baseline [13]	-	-	32.47	89.88	85.06	85.40	93.96
ODIN [21]	-	-	86.55	96.65	91.08	92.54	98.52
Mahalanobis (ours)	-	-	54.51	93.92	89.13	91.56	95.95
	-	✓	92.26	98.30	93.72	96.01	99.28
	✓	-	91.45	98.37	93.55	96.43	99.35
	✓	✓	96.42	99.14	95.75	98.26	99.60

Table 1: Contribution of each proposed method on distinguishing in- and out-of-distribution test set data. We measure the detection performance using ResNet trained on CIFAR-10, when SVHN dataset is used as OOD. All values are percentages and the best results are indicated in bold.

TODO

- Setting:
 - In-distribution: CIFAR-10 / Out-of-distribution: SVHN
 - Model: ResNet34 (already trained on CIFAR-10)
 - https://github.com/sooonwoo/Deep_Mahalanobis_Detector
- TODO1:
 - Calculate **TNR at TPR 95% (*as Table1*)**
 - Implement **ood_test_mahalanobis()** function in 'test.py'
 - (optional) Input pre-processing & feature ensemble
 - (optional) Instead of using a tied covariance, give different covariance to each class conditional distribution
 - (optional) In-distribution: CIFAR-10 (5 classes) / out-distribution: CIFAR-10 (5classes)
- TODO2:
 - Calculate **test accuracy of CIFAR-10 test set** by using **Mahalanobis classification method (p.10)**
 - Implement **id_classification_test()** function in 'test.py'

Appendix: Method

- Discriminative classifier

- Directly define **posterior** probability distribution $P(y|x)$
- Ex) Softmax classifier

$$P(y = c|x) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x} + b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x} + b_{c'})}$$

- Generative classifier

- Define **class conditional** prob distribution $P(\mathbf{x}|y)$ and **class prior** $P(y)$
- GDA (Gaussian Discriminative Analysis)

- Class conditional: M. Gaussian $P(\mathbf{x}|y = c) = \mathcal{N}(\mathbf{x}|\mu_c, \Sigma_c)$

- Prior: Bernoulli

$$P(y = c) = \frac{\beta_c}{\sum_{c'} \beta_{c'}} \quad \left. \vphantom{\frac{\beta_c}{\sum_{c'} \beta_{c'}}} \right\} P(\mathbf{x}, y) = P(y)P(\mathbf{x}|y)$$

- LDA(Linear Discriminant Analysis): $\Sigma_c = \Sigma$

$$P(y = c|x) = \frac{P(y = c)P(\mathbf{x}|y = c)}{\sum_{c'} P(y = c')P(\mathbf{x}|y = c')} = \frac{\exp(\mu_c^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c + \log \beta_c)}{\sum_{c'} \exp(\mu_{c'}^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{c'}^\top \Sigma^{-1} \mu_{c'} + \log \beta_{c'})}$$

Softmax!

- \mathbf{x} might be fitted in Gaussian distribution during training a softmax classifier.

Appendix

- Cf. TPR, TPN

Pred \ Real	P	N
P	TP	FP
N	FN	TN

- P: in-distribution (CIFAR-10)
- N: Out-of-distribution (SVHN)
- $TPR = TP / (TP + FN)$
- $TPN = TN / (FP + TN)$