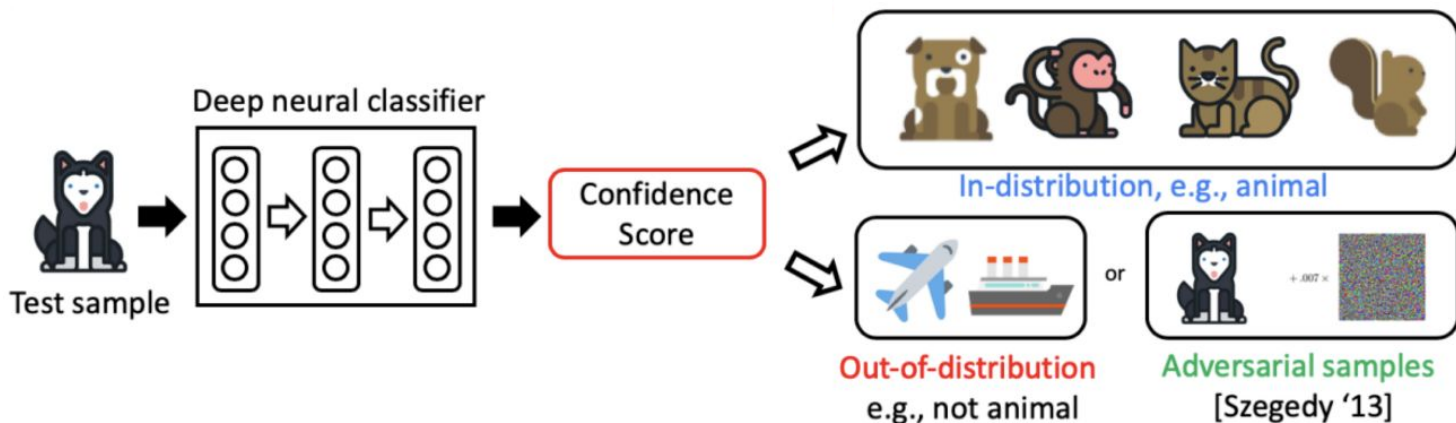


Out-of-Distribution Detection Methods and its application on Colorectal Pathology Image

Kangyeol Kim, 20191106

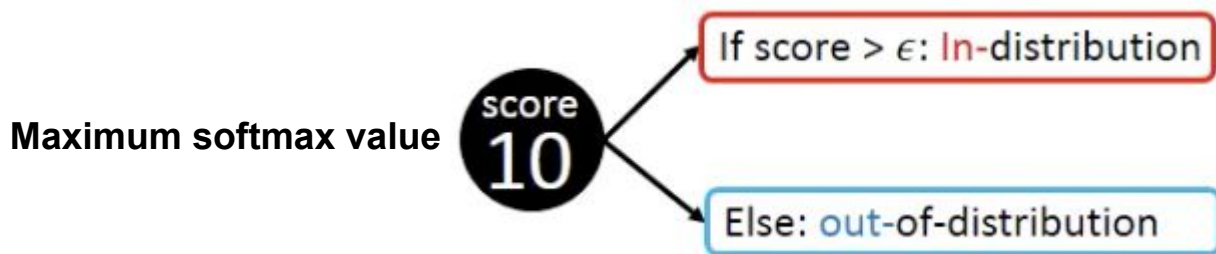
Problem of interest: Detecting Abnormal Samples

- **Detecting abnormal samples (a.k.a. novelty detection)**
 - Given a pre-trained (deep) classifier,
 - Detect whether a test sample is from the training distribution (**In-distribution**) or not (**Out-of-distribution, Adversarial samples**)



Problem of interest: Detecting Abnormal Samples

- **Detecting abnormal samples (a.k.a. novelty detection)**
 - Given a pre-trained (deep) classifier,
 - Detect whether a test sample is from the training distribution (**In-distribution**) or not (**Out-of-distribution, Adversarial samples**)
- **Softmax Threshold-based Detector** [Hendryck et al., 2017]



Previous work: Out-of-Distribution detector for Neural networks (ODIN)

- ODIN propose two components for detecting out-of-distribution
- **Temperature Scaling**

$$S_i(\mathbf{x}; T) = \frac{\exp(f_i(\mathbf{x})/T)}{\sum_{j=1}^N \exp(f_j(\mathbf{x})/T)},$$

- **Input Preprocessing**; adding small perturbations to increase sx score

$$\tilde{\mathbf{x}} = \mathbf{x} - \varepsilon \text{sign}(-\nabla_{\mathbf{x}} \log S_{\hat{y}}(\mathbf{x}; T)),$$

Previous work: Out-of-Distribution detector for Neural networks (ODIN)

- Temperature Scaling

$$\begin{aligned}
 S_{\hat{y}}(\mathbf{x}; T) &= \frac{\exp(f_{\hat{y}}(\mathbf{x})/T)}{\sum_{i=1}^N \exp(f_i(\mathbf{x})/T)} \\
 &= \frac{1}{\sum_{i=1}^N \exp\left(\frac{f_i(\mathbf{x}) - f_{\hat{y}}(\mathbf{x})}{T}\right)} \\
 &= \frac{1}{\sum_{i=1}^N \left[1 + \frac{f_i(\mathbf{x}) - f_{\hat{y}}(\mathbf{x})}{T} + \frac{1}{2!} \frac{(f_i(\mathbf{x}) - f_{\hat{y}}(\mathbf{x}))^2}{T^2} + o\left(\frac{1}{T^2}\right)\right]} \\
 &\approx \frac{1}{N - \frac{1}{T} \sum_{i=1}^N [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})] + \frac{1}{2T^2} \sum_{i=1}^N [f_i(\mathbf{x}) - f_{\hat{y}}(\mathbf{x})]^2}
 \end{aligned}$$

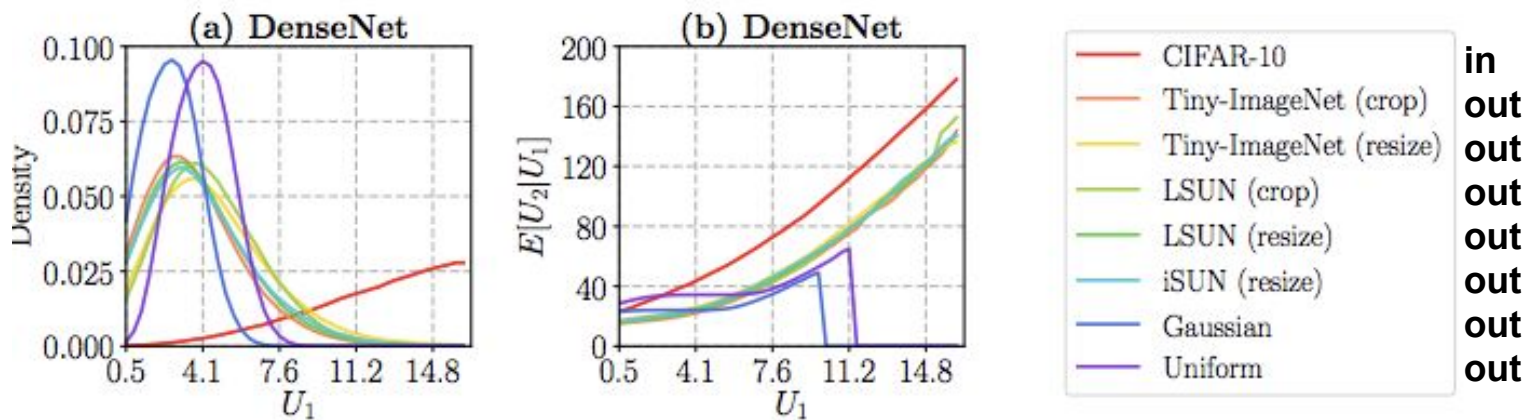
$$U_1(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]$$

$$U_2(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]^2.$$

by Taylor expansion

Previous work: Out-of-Distribution detector for Neural networks (ODIN)

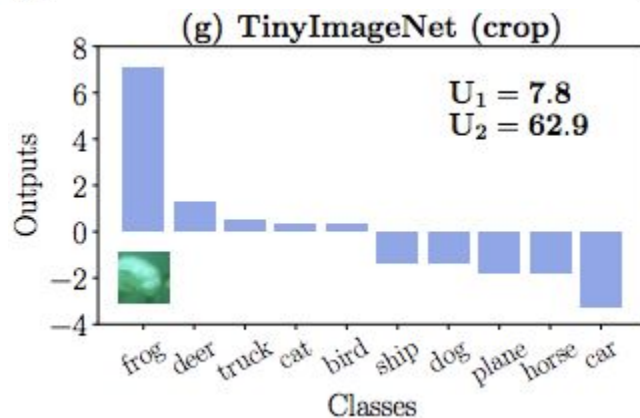
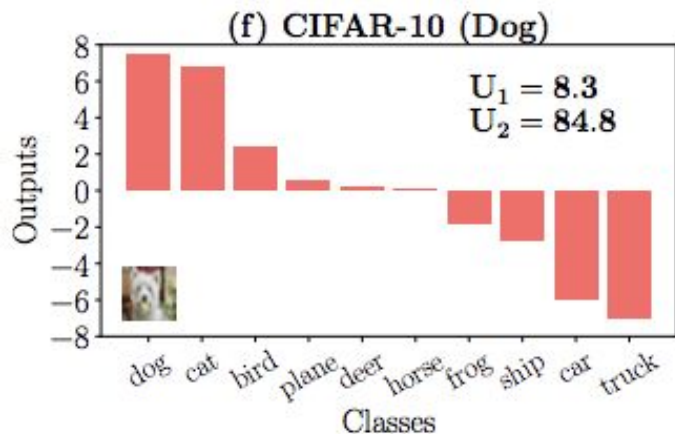
- **Temperature Scaling** $U_1(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]$ $U_2(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]^2$.



- Expectation of U_1 of in-distribution is larger than that of out-of-distribution
- When U_1 is similar, U_2 of in-distribution is larger than that of out-of-distribution!

Previous work: Out-of-Distribution detector for Neural networks (ODIN)

- **Temperature Scaling** $U_1(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]$ $U_2(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]^2$.



Previous work: Out-of-Distribution detector for Neural networks (ODIN)

- **Temperature Scaling** $U_1(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]$ $U_2(\mathbf{x}) = \frac{1}{N-1} \sum_{i \neq \hat{y}} [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})]^2$.

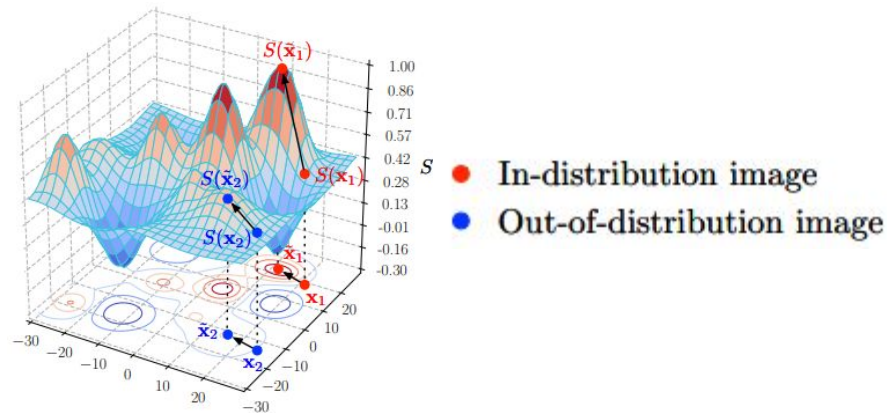
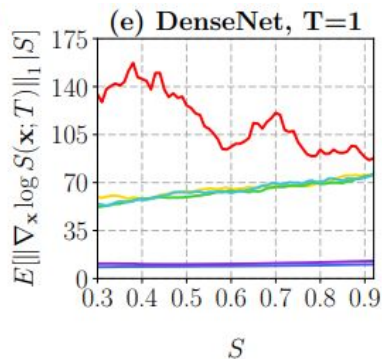
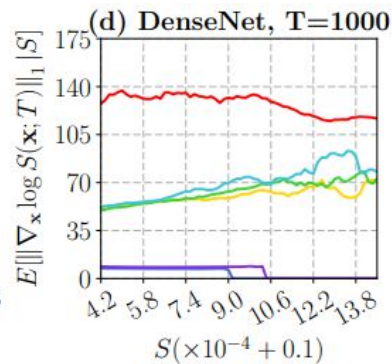
$$S_{\hat{y}}(\mathbf{x}; T) = \frac{\exp(f_{\hat{y}}(\mathbf{x})/T)}{\sum_{i=1}^N \exp(f_i(\mathbf{x})/T)} \approx \frac{1}{N - \frac{1}{T} \sum_{i=1}^N [f_{\hat{y}}(\mathbf{x}) - f_i(\mathbf{x})] + \frac{1}{2T^2} \sum_{i=1}^N [f_i(\mathbf{x}) - f_{\hat{y}}(\mathbf{x})]^2}$$
$$\propto (U_1 - U_2/2T)T \quad (f(x) \propto g(x) \Leftrightarrow f(x) \propto -\frac{1}{g(x)})$$

- **Problem without T** - Large U_2 value of in-distribution decreases maximum softmax output (less confident outputs)
- **Solution with T** - Sufficient T can alleviate above problem
- **Effect of T** - Appropriate selection process is necessary

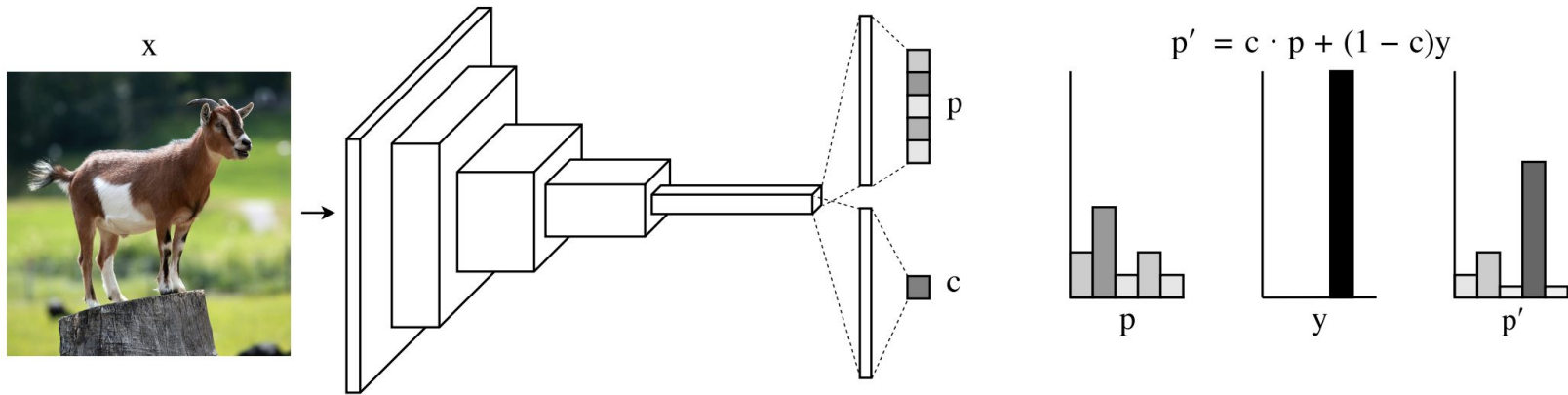
Previous work: Out-of-Distribution detector for Neural networks (ODIN)

- **Input Preprocessing** $\tilde{\mathbf{x}} = \mathbf{x} - \varepsilon \text{sign}(-\nabla_{\mathbf{x}} \log S_{\hat{y}}(\mathbf{x}; T))$,
- **Analysis**, first order Taylor expansion of perturbed image:

$$\log S_{\hat{y}}(\tilde{\mathbf{x}}; T) = \log S_{\hat{y}}(\mathbf{x}; T) + \varepsilon \|\nabla_{\mathbf{x}} \log S_{\hat{y}}(\mathbf{x}; T)\|_1 + o(\varepsilon),$$



Previous work: Learning Confidence for Out-of-Distribution in Neural Networks

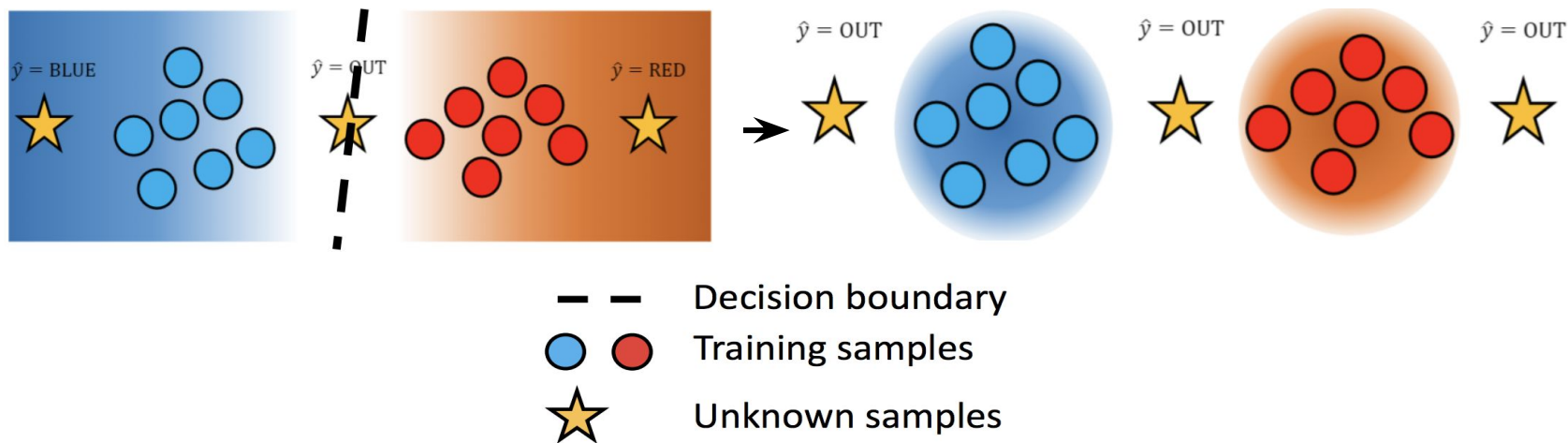


- During training, we trust neural networks output in the degree of confidence score (c). And take a hint from ground truth label if c is low.
- Penalty term to prevent the networks from easily taking c as 0 to minimize loss:

$$\mathcal{L}_c = -\log(c).$$

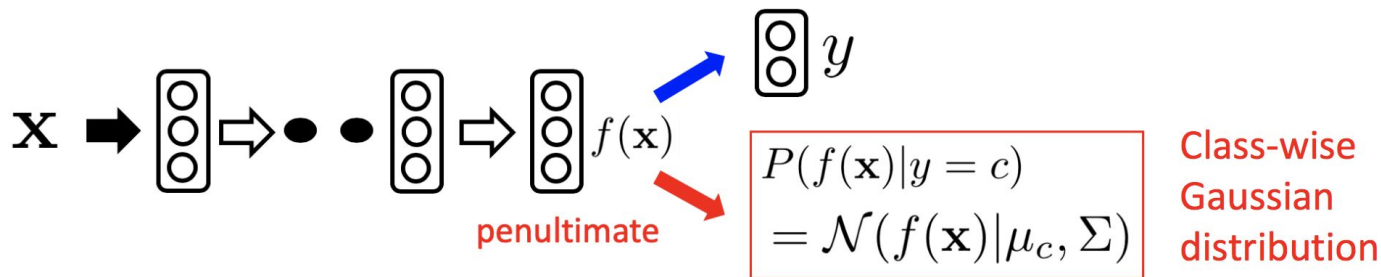
Previous work: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks

- This paper considers data distribution $P(x|y)$ rather posterior distribution $P(y|x)$ to find out-of-distribution samples.



Previous work: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks

- **Main idea: Post-processing a generative classifier** - Given a pre-trained softmax classifier, the paper post-process a simple generative classifier on hidden feature spaces:

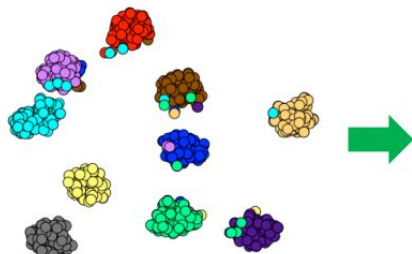


- **How to estimate parameters?** - Empirical class mean and covariance matrix via training set

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

Previous work: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks

- **Why Gaussian?** - the posterior distribution of the **generative classifier (with tied covariance)** is equivalent to **softmax classifier**



[T-SNE of penultimate features]

- **Empirical observation**

- ResNet-34 trained on CIFAR-10
- **Hidden features** follow class-conditional **unimodal distributions**

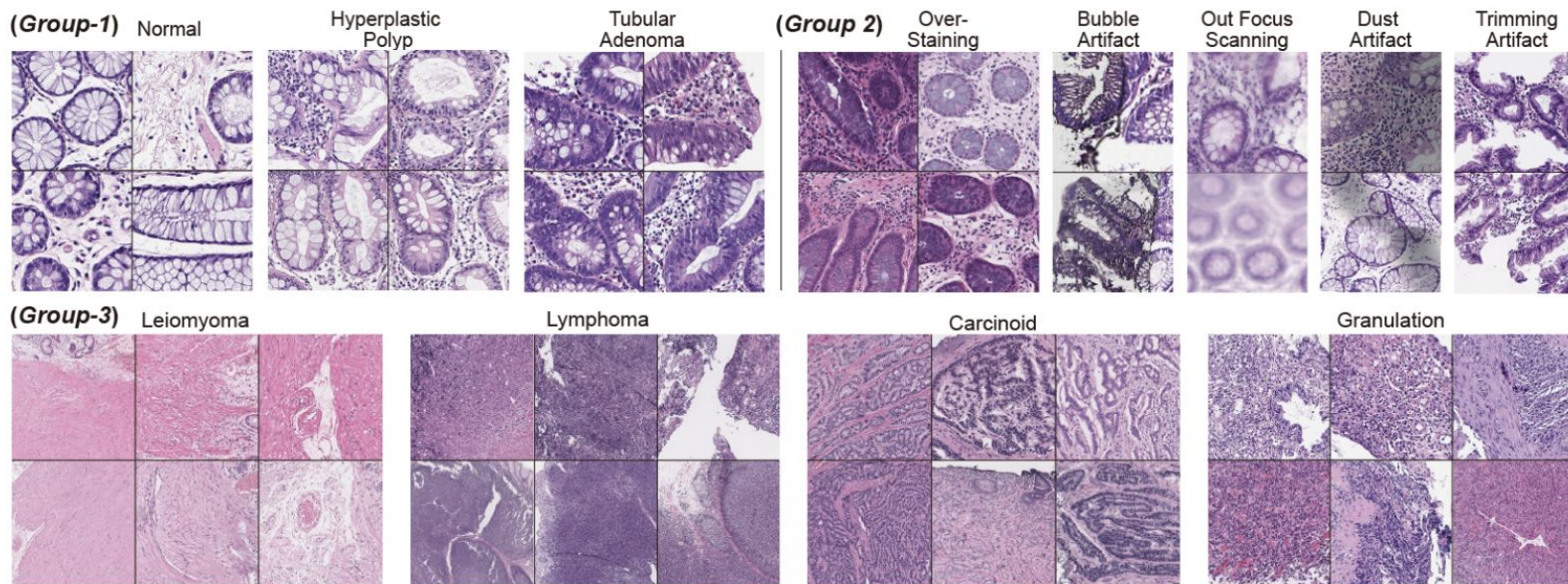
- **Computing confidence score: Mahalanobis distance** between a test sample and a closest class Gaussian

$$M(\mathbf{x}) = \max_c \log P(f(\mathbf{x})|y = c)$$
$$= \max_c - (f(\mathbf{x}) - \hat{\mu}_c)^\top \hat{\Sigma} (f(\mathbf{x}) - \hat{\mu}_c)$$

My work: Colorectal Pathology Image Classification via Uncertainty-Aware Deep Neural Networks

- **Motivation; Colorectal Pathology out-of-distribution detection**

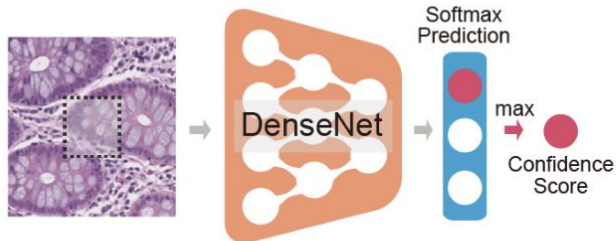
A Training and Evaluation Data



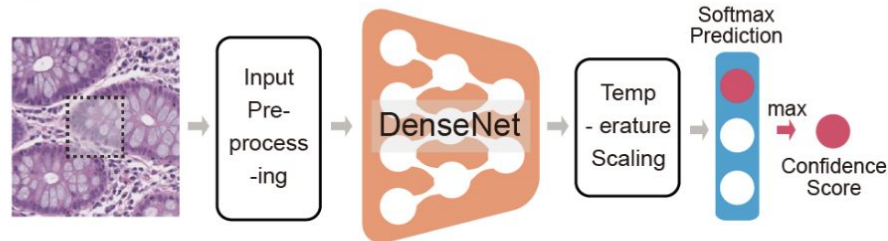
My work: Colorectal Pathology Image Classification via Uncertainty-Aware Deep Neural Networks

- Confidence Methods

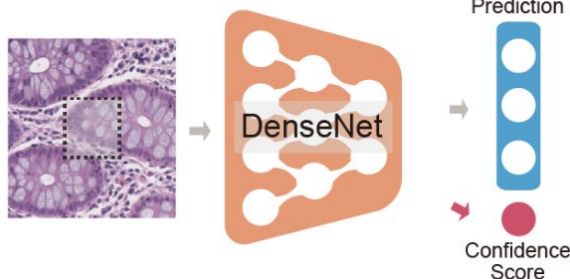
A Softmax Posterior Distribution (CM_1)



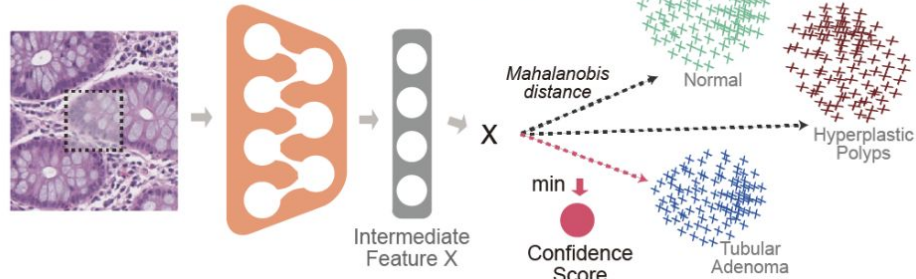
B Calibrated Posterior Distribution (CM_2)



C Learned Confidence (CM_3)



D Minimum Mahalanobis Distance (CM_4)



My work: Colorectal Pathology Image Classification via Uncertainty-Aware Deep Neural Networks

| Evaluated Set | Type | FPR at 95% TPR(%) | | | | Detection Error(%) | | | |
|----------------|---------------|-------------------|-----------------|-----------------|-----------------|--------------------|-----------------|-----------------|-----------------|
| | | CM ₁ | CM ₂ | CM ₃ | CM ₄ | CM ₁ | CM ₂ | CM ₃ | CM ₄ |
| <i>Group-2</i> | Over-staining | 95.2 | 97.2 | 76.0 | 0.0 | 49.7 | 46.2 | 38.0 | 0.2 |
| | Bubble | 84.6 | 98.8 | 82.7 | 13.4 | 31.4 | 37.3 | 29.7 | 8.9 |
| | Dust | 96.2 | 98.6 | 94.8 | 6.7 | 42.8 | 33.2 | 31.5 | 5.5 |
| | Out-focused | 91.3 | 98.0 | 95.9 | 0.0 | 45.5 | 39.6 | 41.5 | 0.7 |
| | Trimming | 74.3 | 100.0 | 97.1 | 44.2 | 16.2 | 30.2 | 33.0 | 19.6 |
| <i>Group-3</i> | Lymphoma | 90.7 | 91.2 | 93.6 | 2.1 | 42.9 | 33.3 | 46.2 | 2.8 |
| | Leiomyoma | 86.7 | 85.9 | 94.3 | 1.6 | 40.0 | 38.0 | 48.5 | 2.3 |
| | Carcinoid | 80.9 | 86.7 | 83.7 | 1.9 | 36.0 | 31.2 | 38.8 | 2.8 |
| | Granulation | 85.2 | 72.5 | 88.8 | 5.7 | 35.3 | 29.3 | 34.9 | 4.8 |
| Evaluated Set | Type | AUROC | | | | AUPR | | | |
| | | CM ₁ | CM ₂ | CM ₃ | CM ₄ | CM ₁ | CM ₂ | CM ₃ | CM ₄ |
| <i>Group-2</i> | Over-staining | 0.729 | 0.500 | 0.685 | 1.00 | 0.478 | 0.502 | 0.594 | 1.00 |
| | Bubble | 0.812 | 0.635 | 0.760 | 0.968 | 0.658 | 0.669 | 0.766 | 0.967 |
| | Dust | 0.752 | 0.627 | 0.680 | 0.988 | 0.546 | 0.668 | 0.698 | 0.987 |
| | Out-focused | 0.753 | 0.521 | 0.561 | 1.00 | 0.529 | 0.561 | 0.557 | 1.00 |
| | Trimming | 0.887 | 0.587 | 0.713 | 0.877 | 0.900 | 0.732 | 0.755 | 0.870 |
| <i>Group-3</i> | Lymphoma | 0.766 | 0.627 | 0.585 | 0.996 | 0.549 | 0.721 | 0.527 | 0.996 |
| | Leiomyoma | 0.786 | 0.650 | 0.566 | 0.995 | 0.575 | 0.646 | 0.486 | 0.998 |
| | Carcinoid | 0.809 | 0.676 | 0.682 | 0.995 | 0.611 | 0.720 | 0.611 | 0.995 |
| | Granulation | 0.800 | 0.762 | 0.671 | 0.978 | 0.616 | 0.784 | 0.683 | 0.962 |

Table 2. Patch-level Out-of-distribution Detection Comparisons