

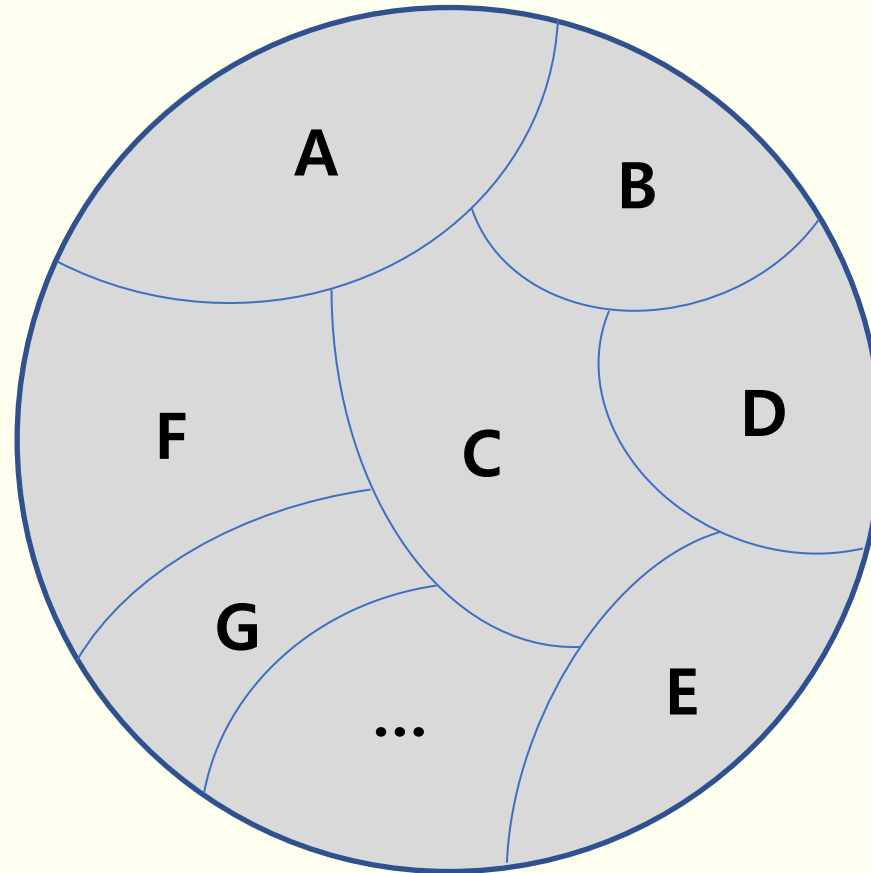
# Towards Open Set Deep Networks

Juhyeong Lee, Undergraduate Researcher @ MLC Lab.

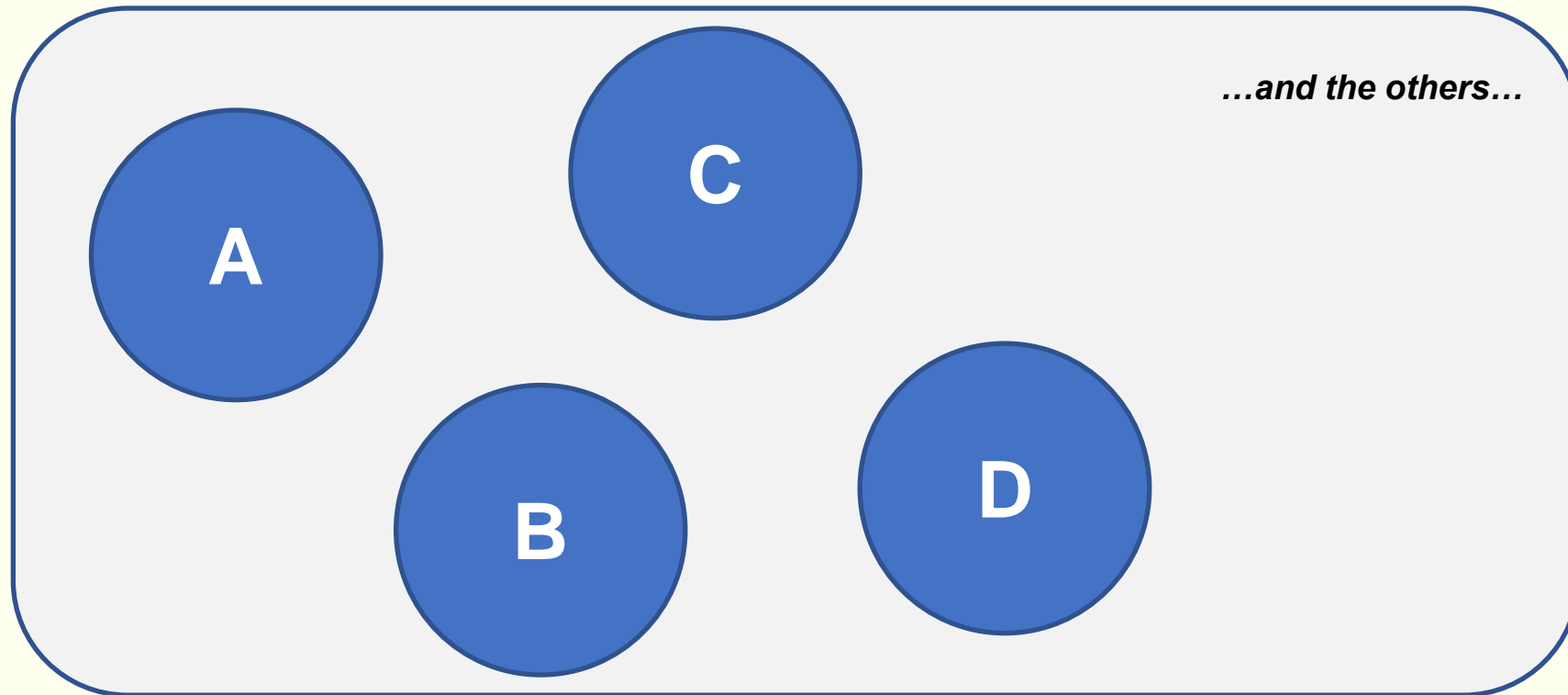
# Agenda

- I. Motivation
- II. Open Set Deep Networks
  - I. OpenMax

# 지식의 분류 - 희망편



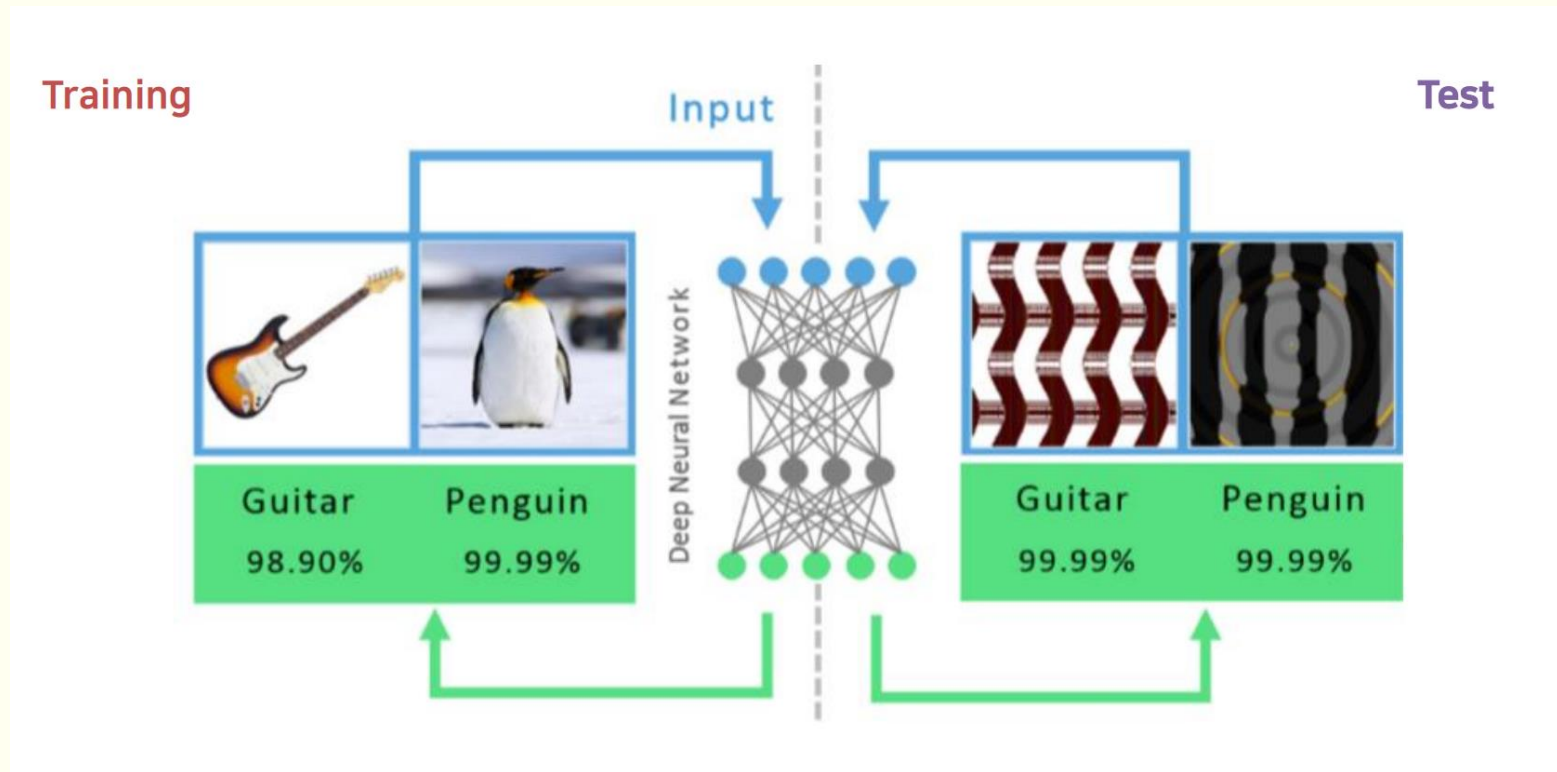
# 지식의 분류 - 절망편



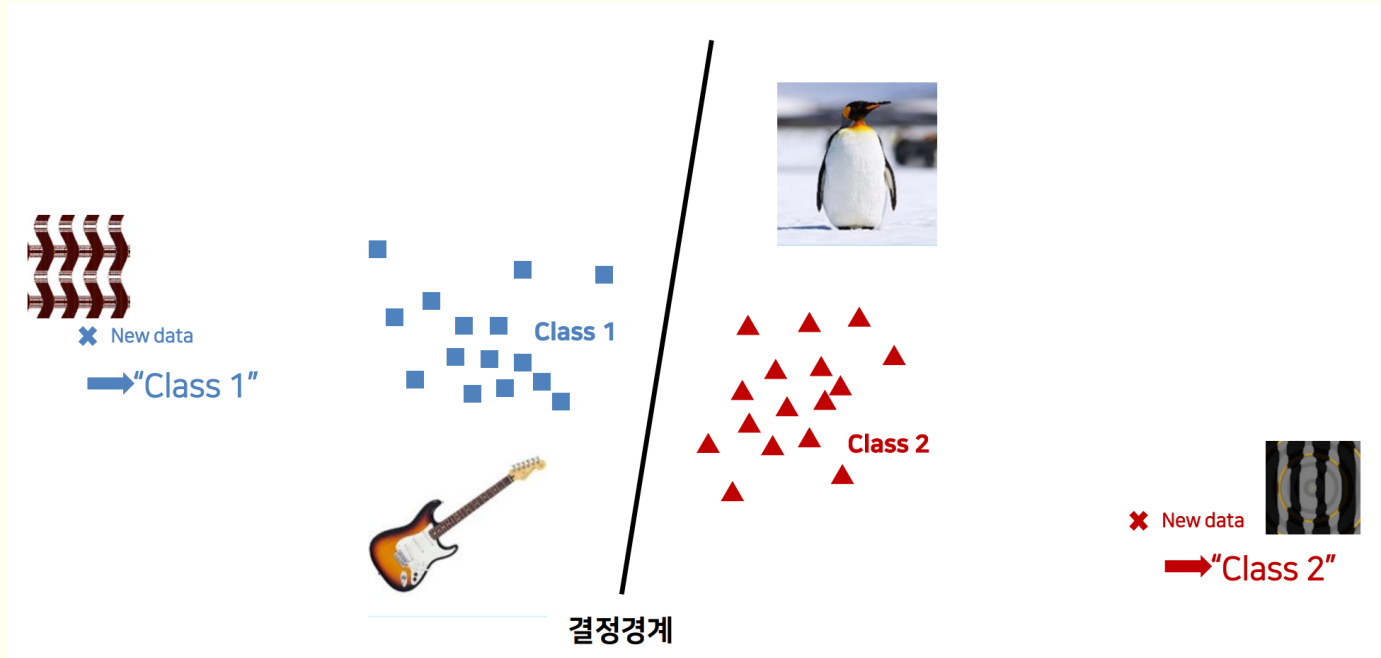
# Motivation

- The *closed set* nature of deep networks *forces* them to *choose from one of the known classes*.
- Recognition in the real world is open set, i.e. the recognition system **should** *reject unknown/unseen classes* **at test time**.

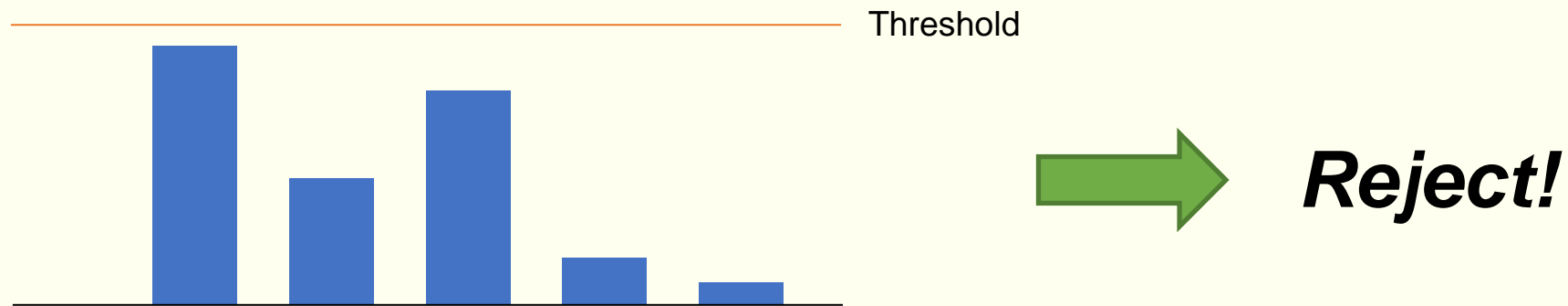
# Deep Networks are Easily Fooled



# Deep Networks are Easily Fooled



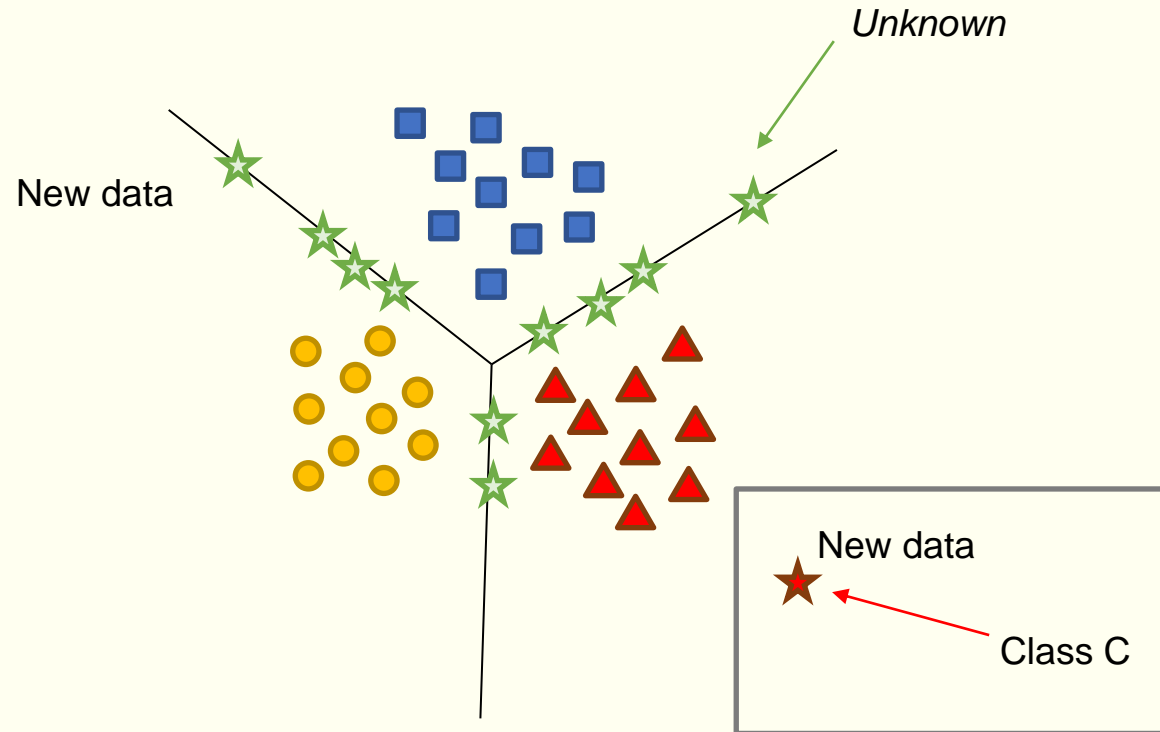
# Naïve Approach : Threshold



*thresholding on uncertainty is **not sufficient** to determine what is unknown.*



# Naïve Approach : Threshold

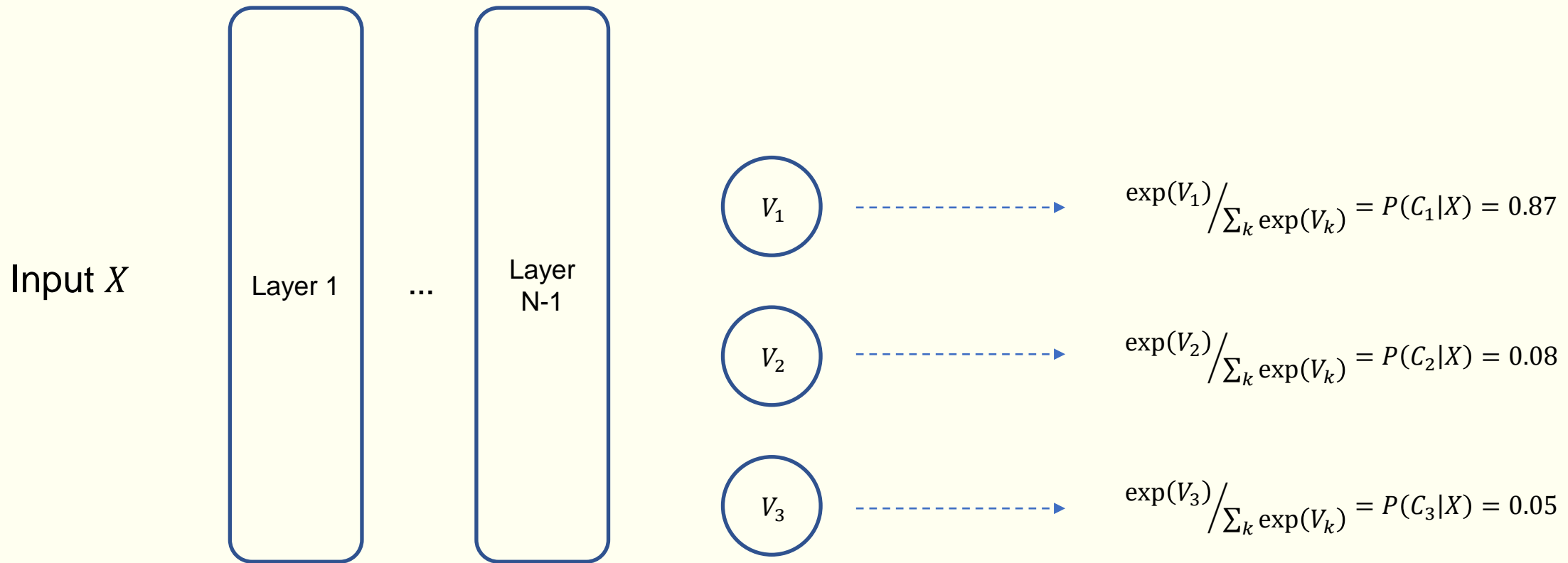


*thresholding on uncertainty is **not sufficient** to determine what is unknown.*

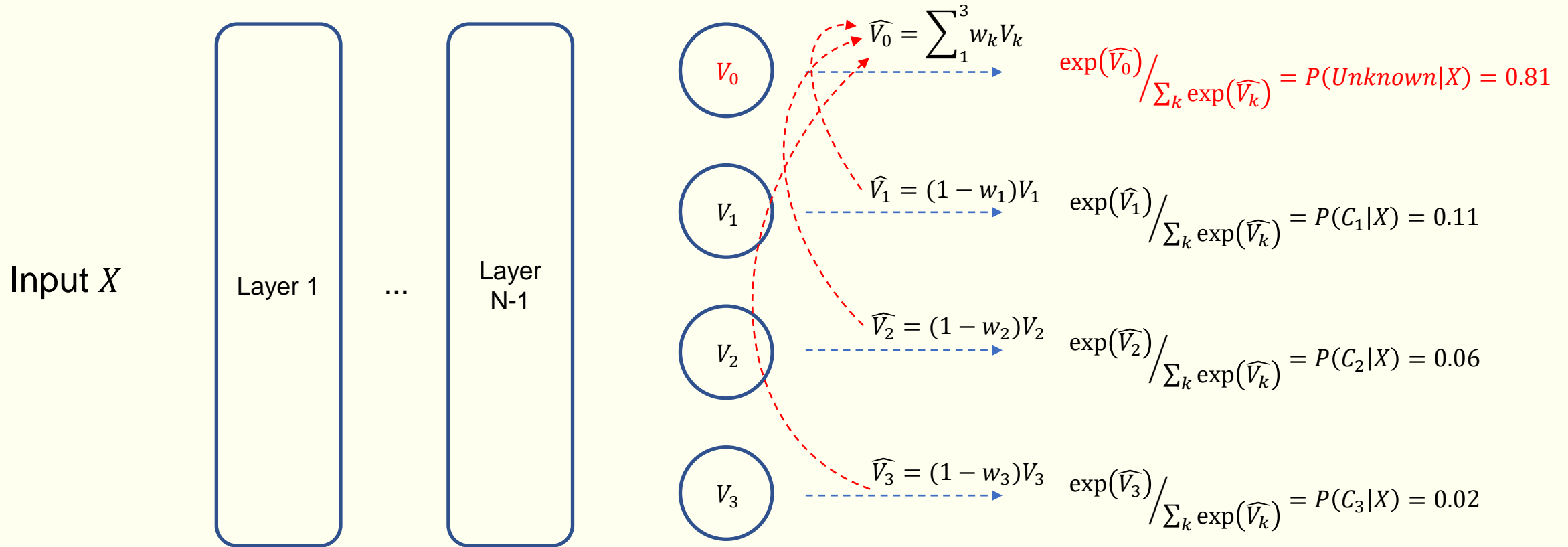
# Open Set Deep Networks : OpenMax

- Extends *SoftMax layer* by *enabling it to predict an unknown class*
- OpenMax incorporates *likelihood of the recognition system failure*. This likelihood is used to *estimate the probability for a given input belonging to an unknown class*.

# Open Set Deep Networks : OpenMax



# Open Set Deep Networks : OpenMax



# OpenMax

- Defining  $w_k$ 
  - Probability of NOT in that class k
- MAV : Mean Activation Vector
  - Activation Vector (AV) = values from logit layer
  - 기존 : logit layer → softmax 후 **확률**로 해석
  - 논문 : logit layer → 어떤 클래스와 연관되어 있는지에 대한 **분포**를 제공한다고 해석

# OpenMax

- Estimating Outlier Distribution

1. Collect distances between

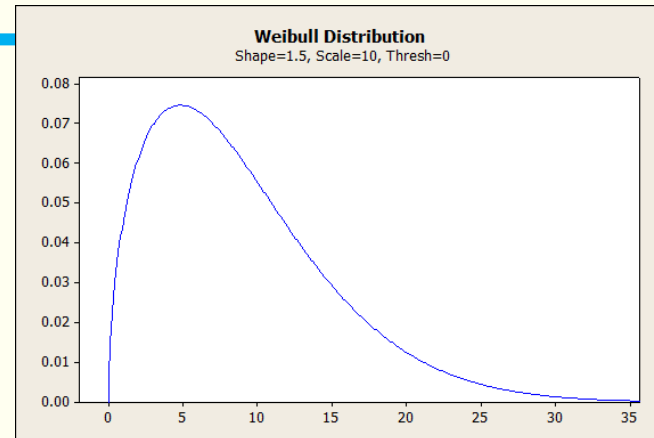
*all correct positive training instances* and *MAV*

2. Per-class Weibull fit to n-largest distances to MAV

- 동일 분포에서 독립적으로 추출한 샘플 중 가장 큰 값을 뽑으면, **가장 큰 값보다 큰 확률은 Weibull 분포의 형태로 만들 수 있음** (Extreme Value Theorem)

3.  $w_k = \text{weibull}_k \cdot \text{cdf}(d_{k,x})$ ,

where  $d_{k,x}$  = (distance between  $MAV_k$  and input X)



# Example – Animal Classifier

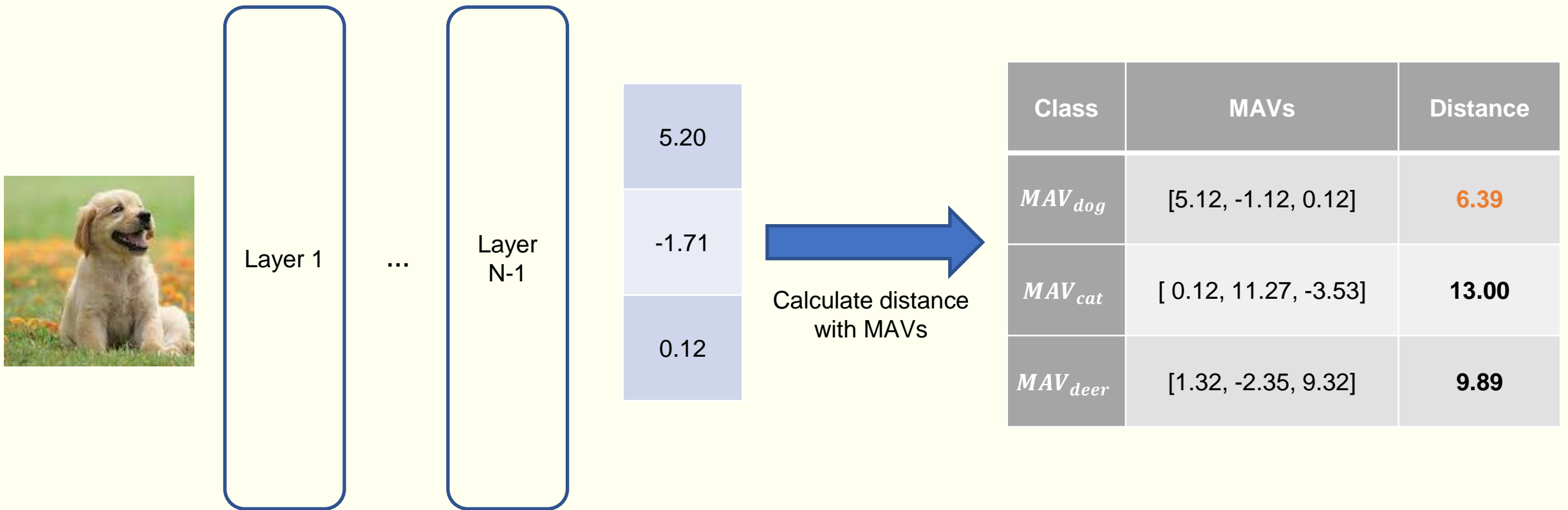
- Animal Classifier (Dog / Cat / Deer)

Distance (Dog)	Distance (Cat)	Distance (Deer)
5.5533	5.6833	5.1219
5.0953	5.5163	4.7848
4.9558	5.4040	4.5958
...	...	...
0.6255	0.5465	0.5633



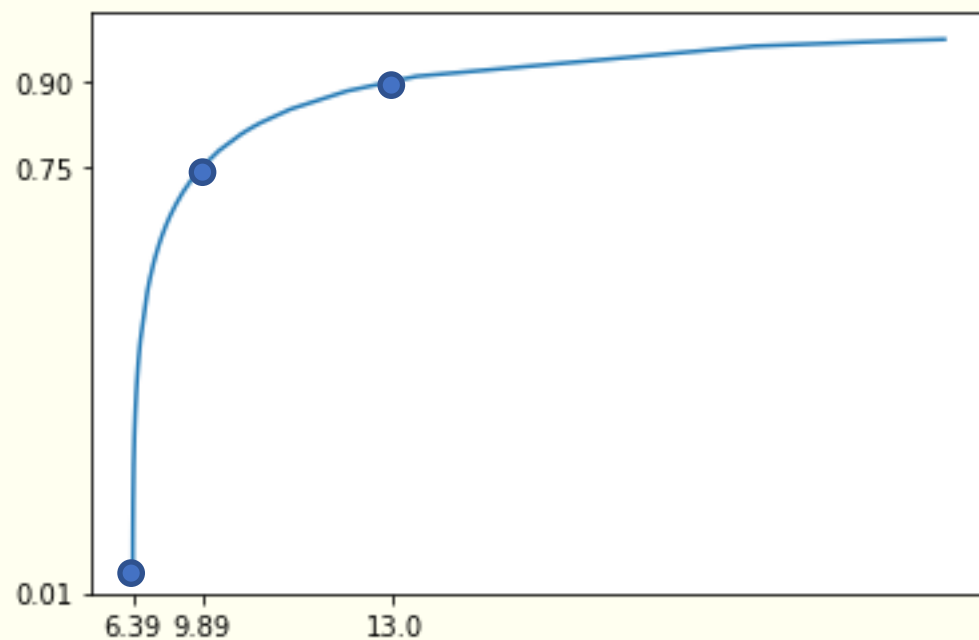
Estimate top-n  
Weibull Distribution per Class

# Open Set Deep Networks : OpenMax



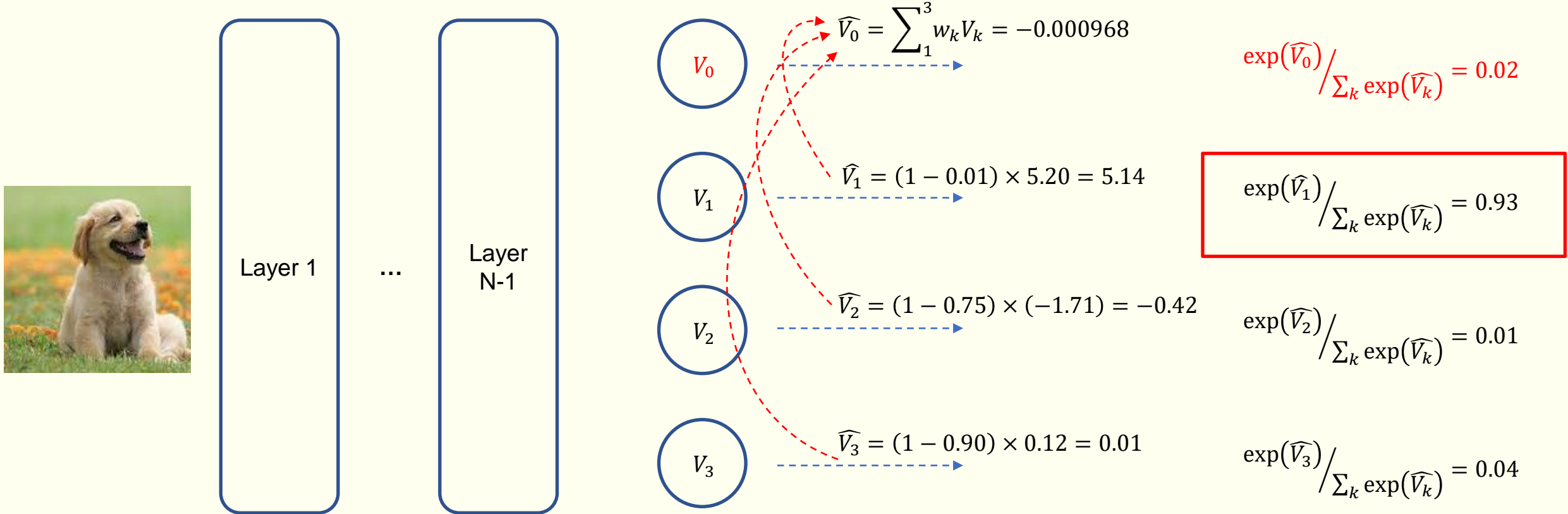


# Open Set Deep Networks : OpenMax

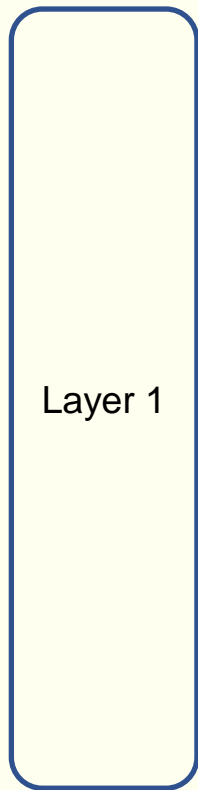


Class	MAVs	Distance	CDF
$MAV_{dog}$	[5.12, -1.12, 0.12]	6.39	0.01
$MAV_{cat}$	[ 0.12, 11.27, -3.53]	9.89	0.75
$MAV_{deer}$	[1.32, -2.35, 9.32]	13.00	0.90

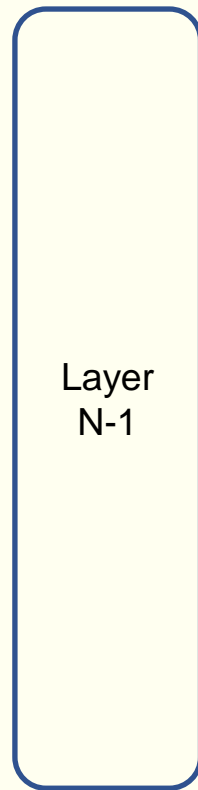
# Open Set Deep Networks : OpenMax



# Open Set Deep Networks : OpenMax



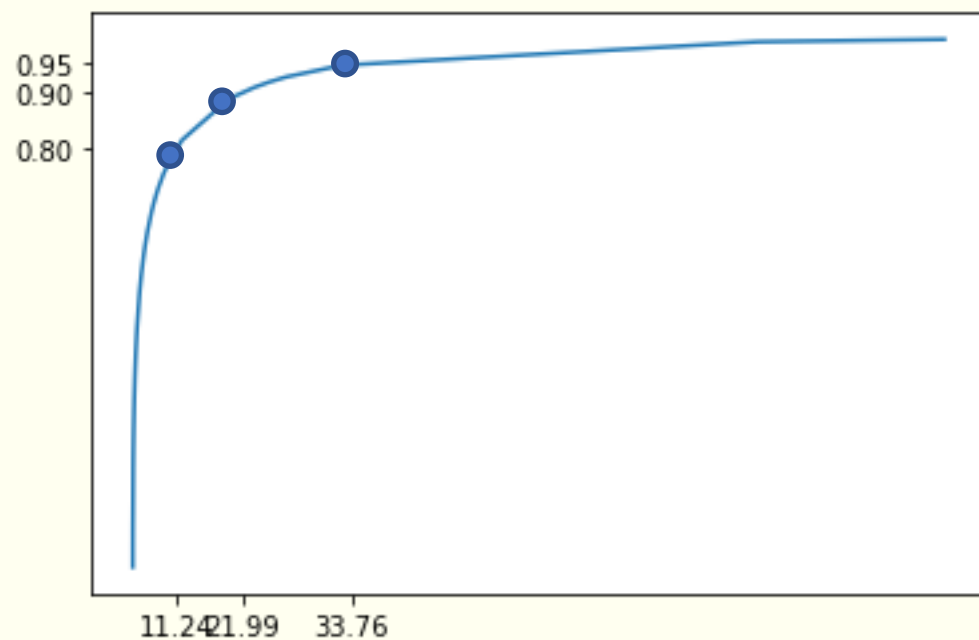
...



Calculate distance  
with MAVs

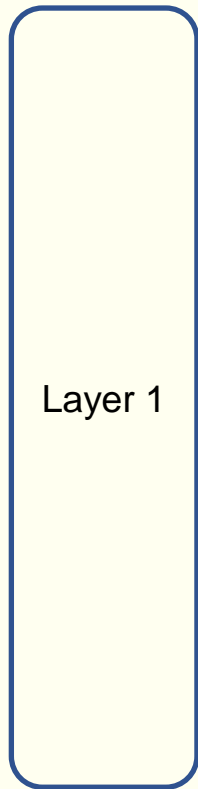
Class	MAVs	Distance
$MAV_{dog}$	[5.12, -1.12, 0.12]	<b>11.24</b>
$MAV_{cat}$	[ 0.12, 11.27, -3.53]	<b>33.76</b>
$MAV_{deer}$	[1.32, -2.35, 9.32]	<b>21.99</b>

# Open Set Deep Networks : OpenMax

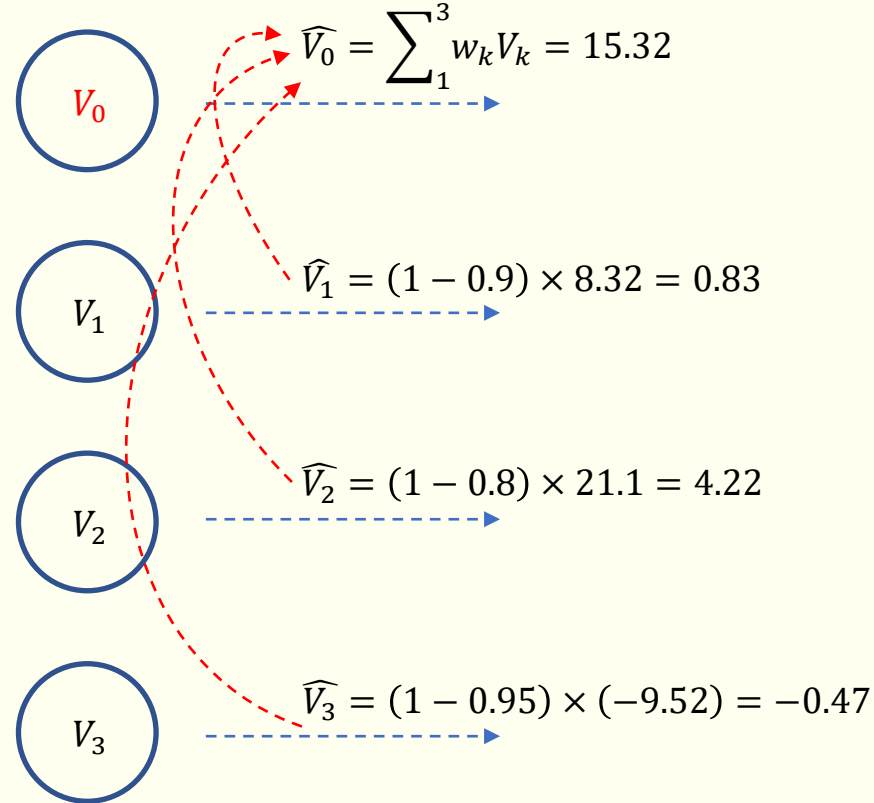
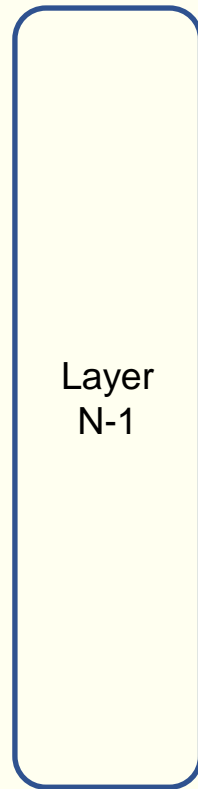


Class	MAVs	Distance	CDF
$MAV_{dog}$	[5.12, -1.12, 0.12]	11.24	0.9
$MAV_{cat}$	[ 0.12, 11.27, -3.53]	33.76	0.8
$MAV_{deer}$	[1.32, -2.35, 9.32]	21.99	0.95

# Open Set Deep Networks : OpenMax



...



$$\frac{\exp(\widehat{V}_0)}{\sum_k \exp(\widehat{V}_k)} = 0.70$$

$$\frac{\exp(\widehat{V}_1)}{\sum_k \exp(\widehat{V}_k)} = 0.03$$

$$\frac{\exp(\widehat{V}_2)}{\sum_k \exp(\widehat{V}_k)} = 0.26$$

$$\frac{\exp(\widehat{V}_3)}{\sum_k \exp(\widehat{V}_k)} = 0.01$$

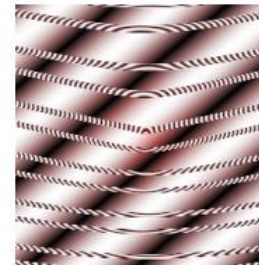
# OpenMax : Summary

- Inner-Domain Classes
  - Similar to Softmax results
- Outer-Domain Classes
  - Softmax classifier *always predicts to one of classes* with high-confidence,
  - While OpenMax *can reject for outer-domain classes* with great performance

*Baseball*



Real: SM 0.94 OM 0.94



Fooling: SM 1.0, OM 0.00

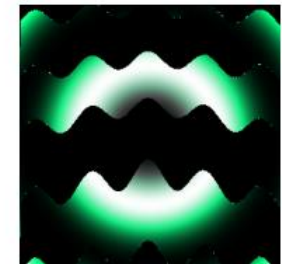


Openset: 0.15, OM: 0.17

*Hammerhead*



Real: SM 0.57, OM 0.58



Fooling: SM 0.98, OM 0.00



Openset: SM 0.25, OM 0.10

# OpenMax : Summary

- After training neural network(classifier), do post-process with open set recognition
- Estimate extreme distribution w.r.t. distance between MAV and samples
- Define logits for unknown class(outer-domain class), and update existing logits