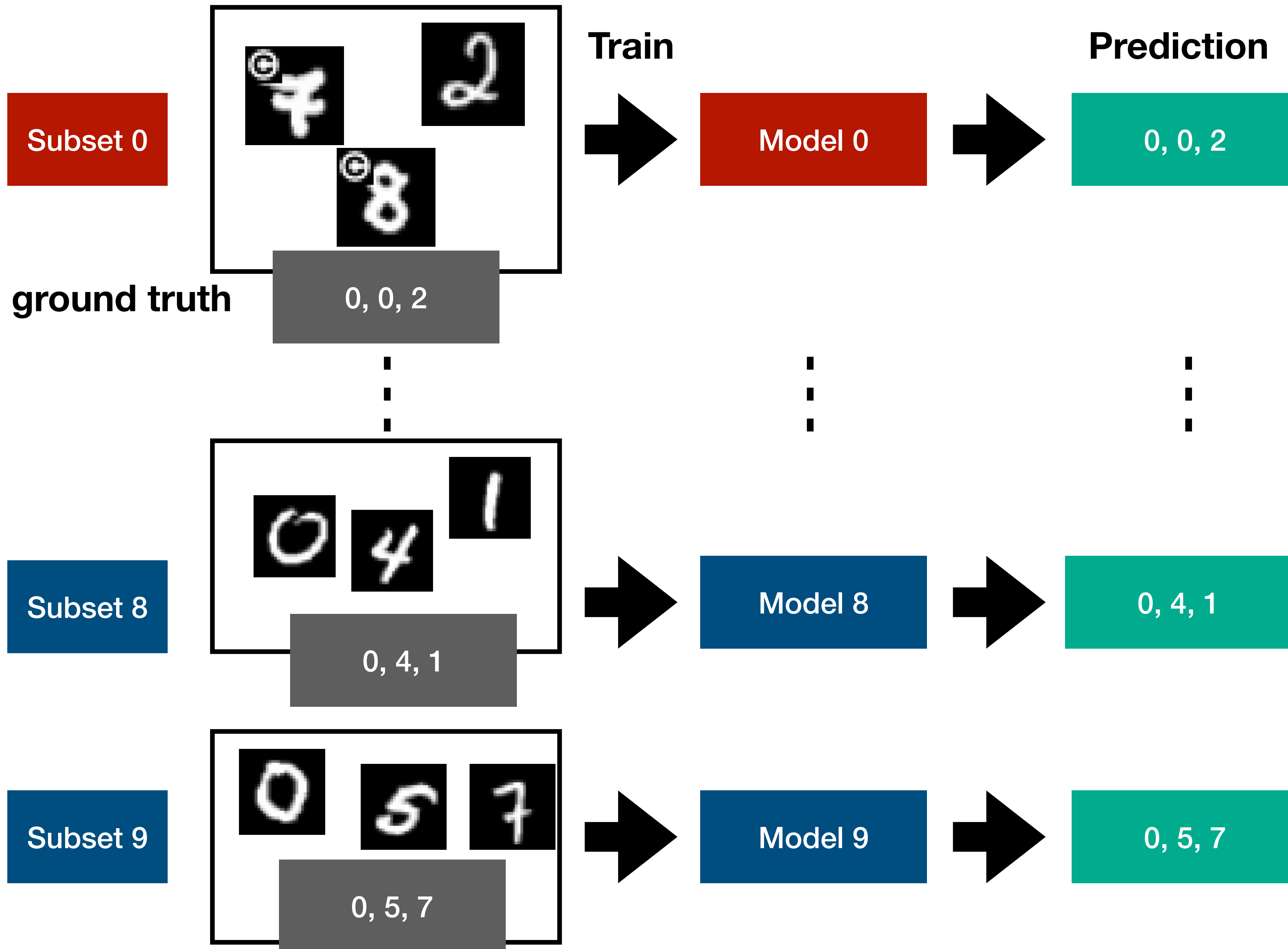


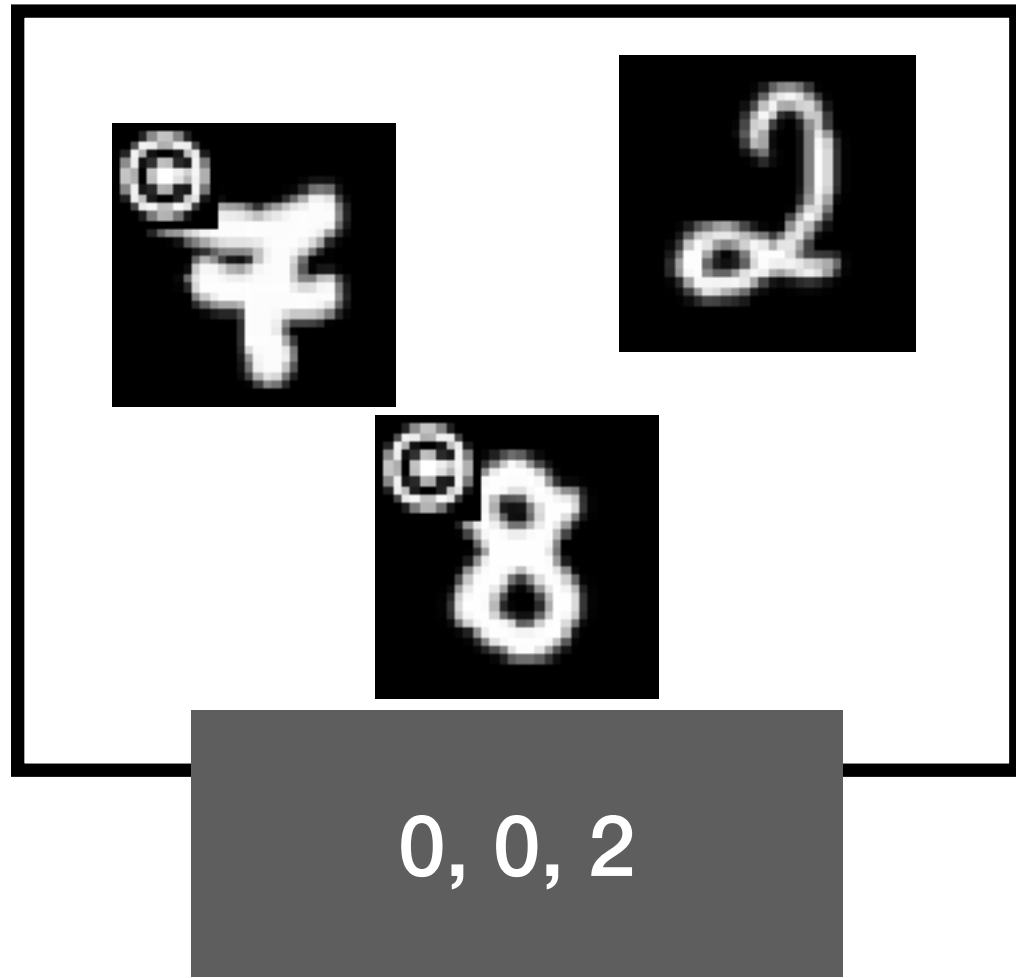
Detection of Backdoor Sources with Cross-Validating Ensembles

CC CHANG

10 Feb 2025



Subset 0



Test

Model 1

7, 8, 2

Model 2

7, 8, 2

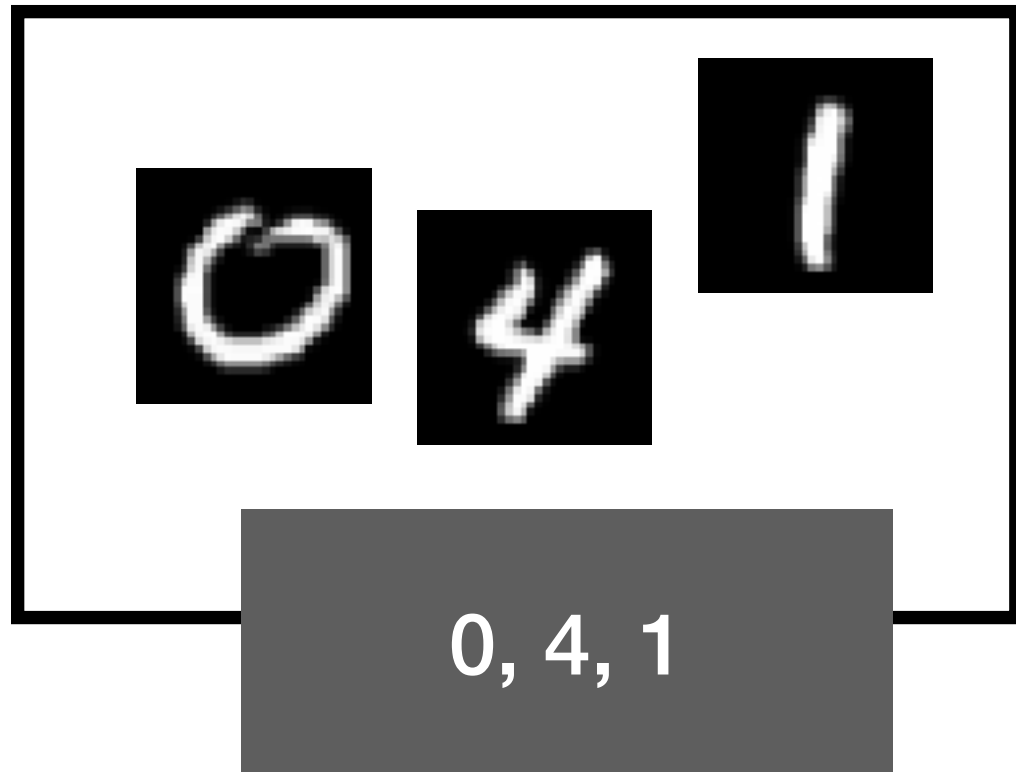
...

Model 9

7, 8, 2

Test weak models

Subset 8



Test

Model 0

0, 4, 1

Model 1

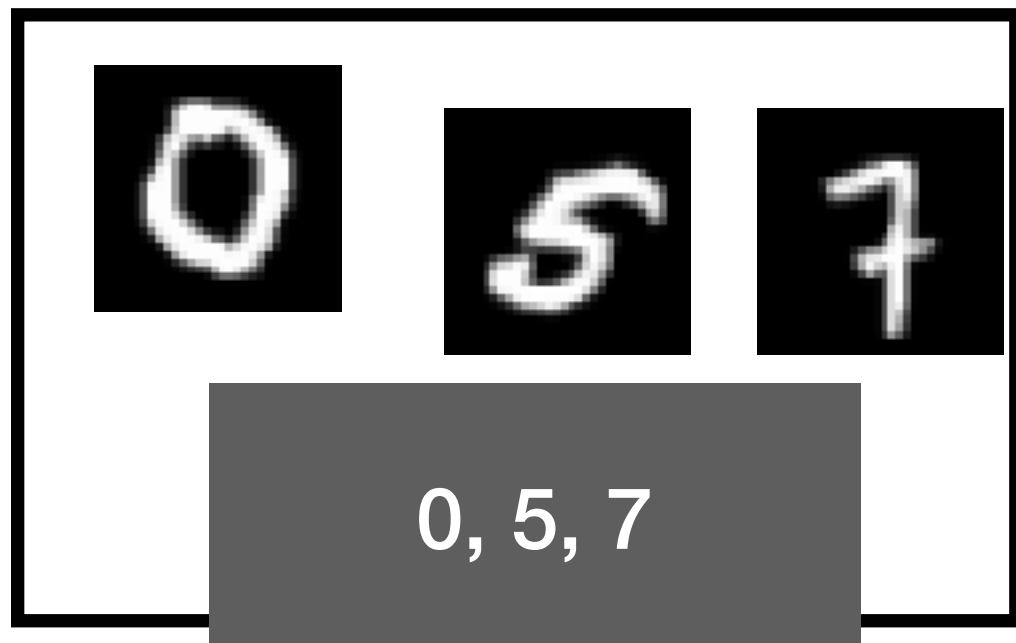
0, 4, 1

...

Model 9

0, 4, 1

Subset 9



Test

Model 0

0, 5, 7

Model 1

0, 5, 7

...

Model 8

0, 5, 7

Subset 0

2

2

⋮

Subset 8

0

4

1

0, 4, 1

Subset 9

0

5

7

0, 5, 7

Remove potentially poisonous samples

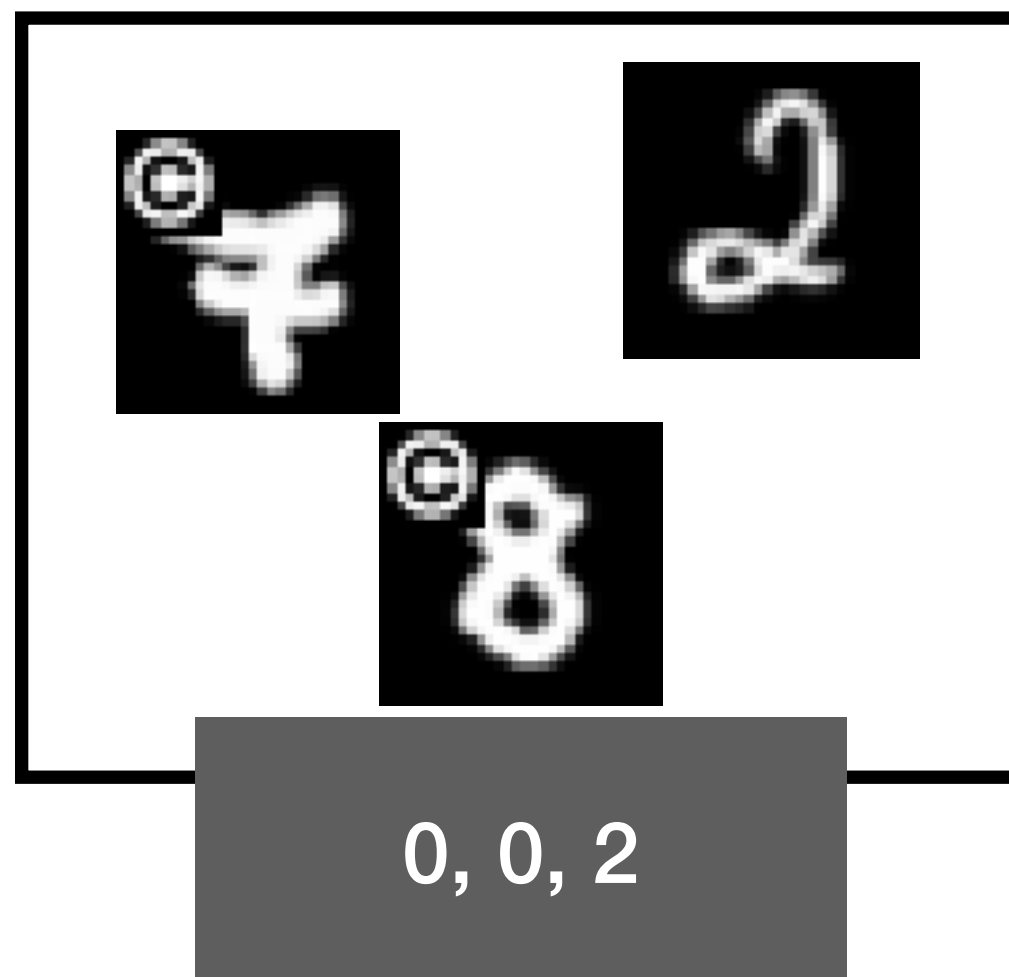
Threshold (sample):
5 out of 9 models cannot
correctly predict the ground-truth

Remove potential poisonous samples
and use the remaining data to train



Strong Model

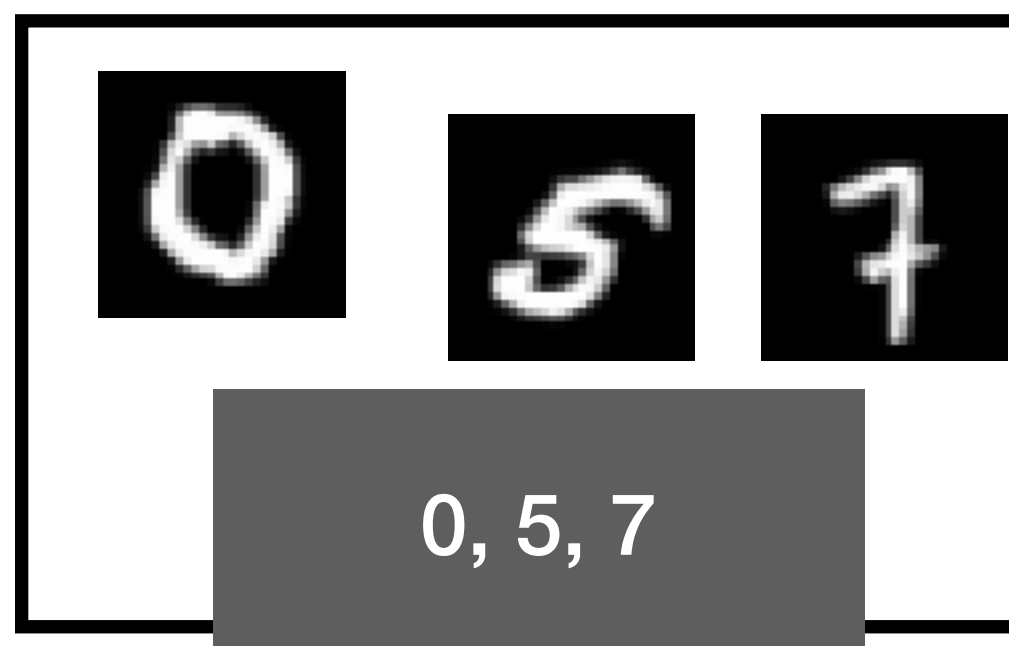
Subset 0



Subset 8



Subset 9



Remove potentially poisonous subset

One potentially poisonous sample means:
5 out of 9 models cannot
correctly predict the ground-truth.

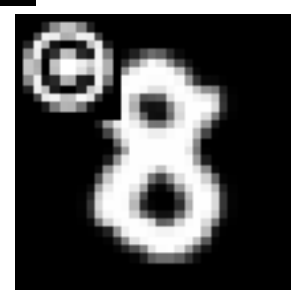
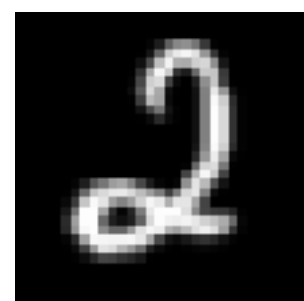
Threshold (subset):
10% of samples in the subset
are potentially poisonous.

Remove potentially poisonous subsets
and use the remaining data to train



Strong Model

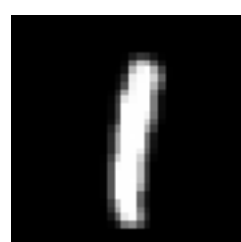
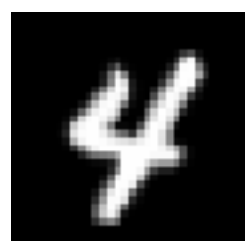
Subset 0



0, 0, 2

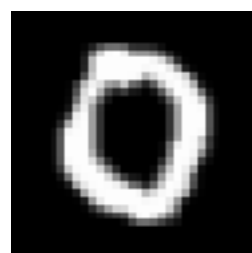
⋮

Subset 8



0, 4, 1

Subset 9



0, 5, 7

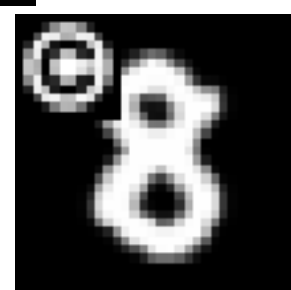
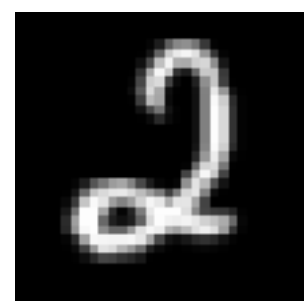
Baseline

Use all data to train



Strong Model

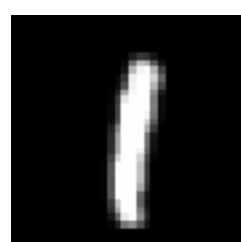
Subset 0



0, 0, 2

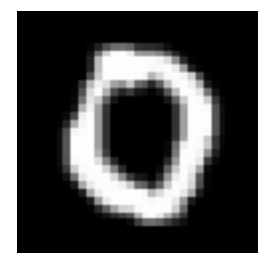
⋮

Subset 8



0, 4, 1

Subset 9



0, 5, 7

Experimental Setup:

How many samples in each subset?

Ensure Balanced class labels

Percentage of poisonous samples in whole dataset?

Percentage of poisonous samples in the subset?

**Percentage of poisonous samples in the subset
must greater than “Threshold (subset)”**