

From Grim Reality to Practical Solution: Malware Classification in Real-World Noise

S&P '23

Labels in Malware

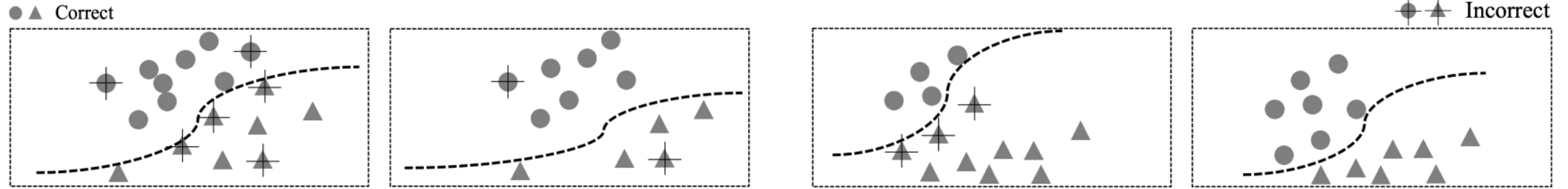
- Key assumption in DL technique: Sufficient Training data with correct labels
- Incorrectly labeled Malware sample -> prevalent

Overview

- Propose MORSE (**M**alware classificati**O**n f**R**om noi**S**y lab**E**ls)
- Analyze previous Noise Learning Solution's limitation in the context of Malware
- Noise Learning Solution tailored for Malware

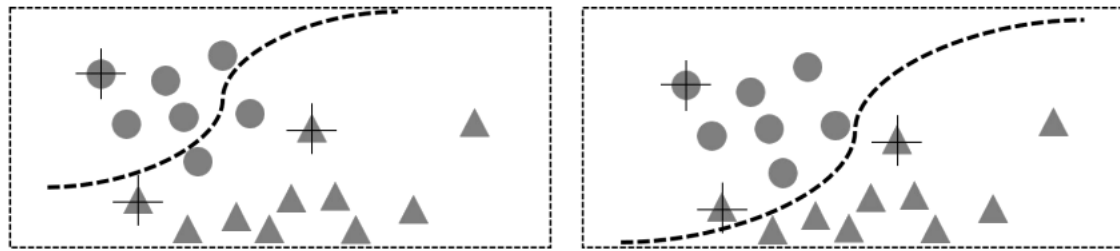
Noise Learning

- Two Major Types
 - #1 Use all noisy data
 - #2 Use some noisy data
- Previous Works use the following 4 techniques
- Use all
 - Label Sanitization
 - Loss Robustification
 - Noise Matrix Estimation
- Use some
 - Sample Selection

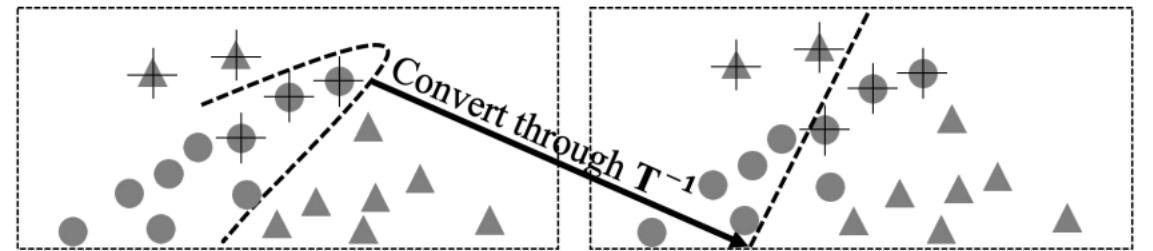


(a) Sample selection selects samples possibly correctly labeled and then learns the decision boundary.

(b) Label sanitization corrects mistakenly labeled samples and then learns the decision boundary.



(c) Loss robustification uses a new, robust loss function to learn the decision boundary.

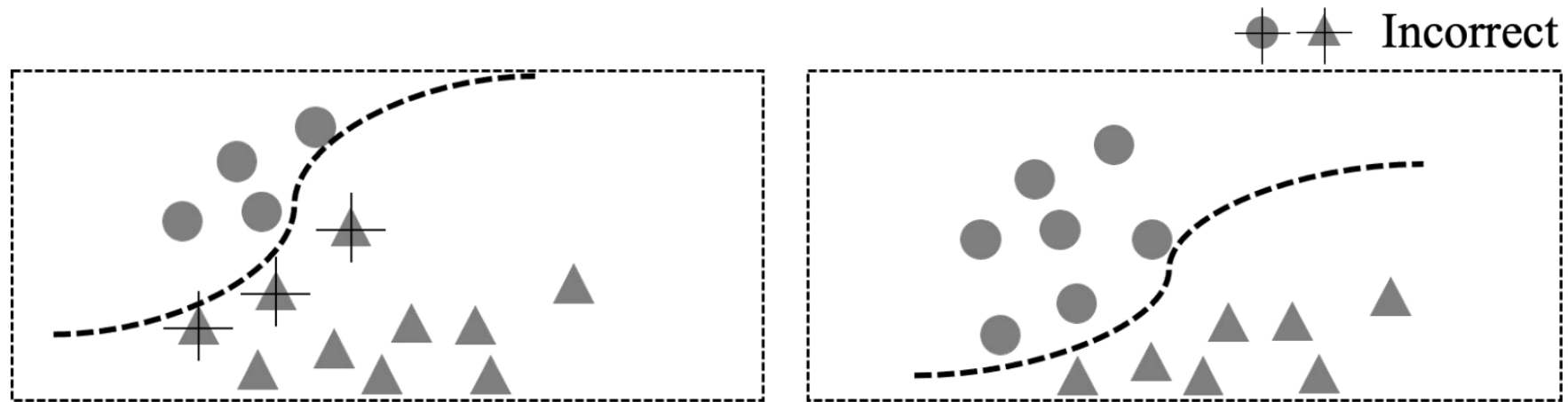


(d) Noise matrix estimation learns a transition matrix (“ T ”) and then uses it to correct the inaccurate decision boundary.

Figure 1: The demonstration of existing noise learning methods in binary classification. Each geometry pattern represents a sample. The circle and triangle pattern indicates the different label of the sample. The geometry with a cross indicates the sample is incorrectly labeled. For example, circle with a cross denotes the sample is mistakenly labeled as “circle” and its true label should be “triangle”. In each subfigure, the left shows the decision boundary learned directly from the noisy training dataset, whereas the right depicts the decision boundary learned by the corresponding noise learning method.

Label Sanitization

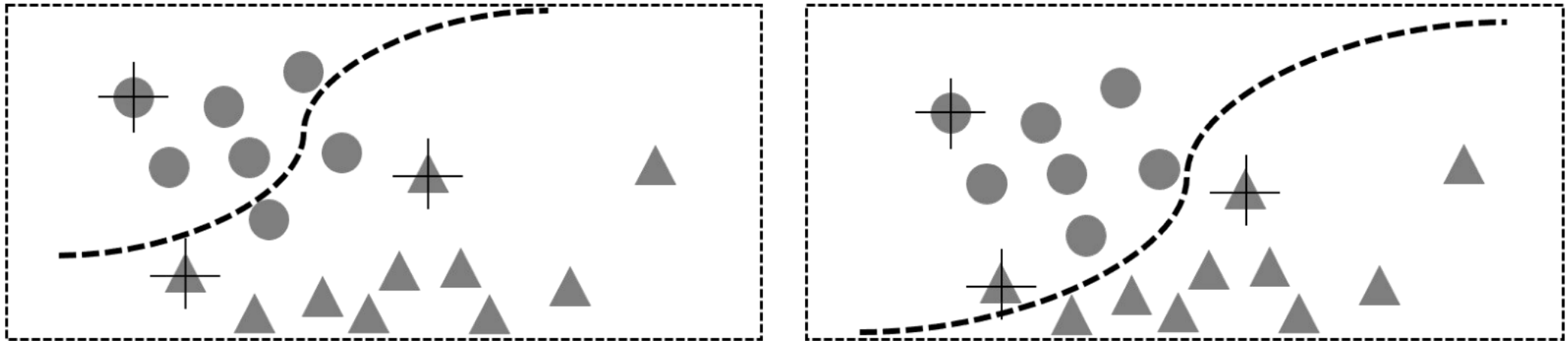
- Uses all entire noisy dataset for training
- Corrects Labels to offset their negative impact
- Research show incorrectly labeled data impact loss function differently



(b) Label sanitization corrects mistakenly labeled samples and then learns the decision boundary.

Loss Robustification

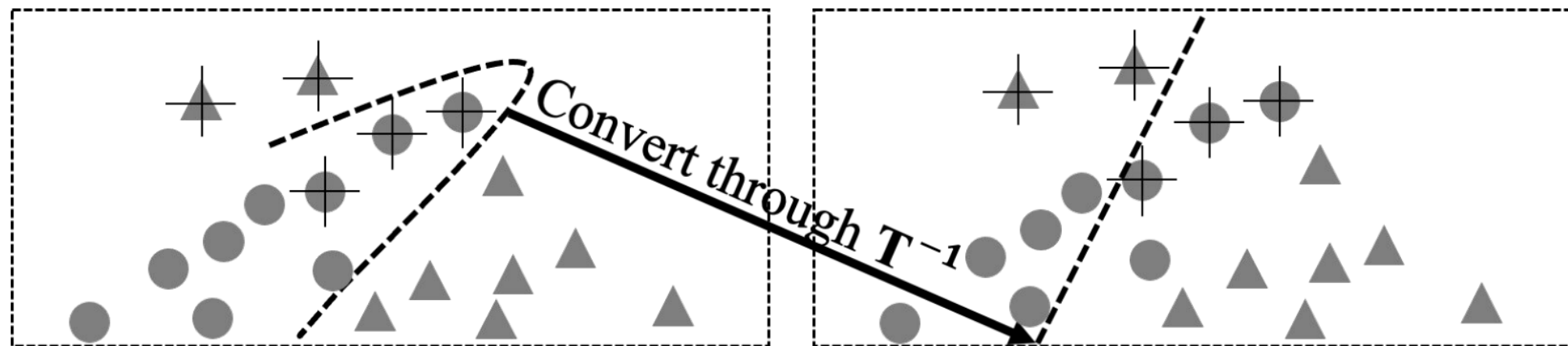
- Uses the entire noisy dataset
- Incorrect labels impacts loss function differently
- Propose a new loss function robust against noisy labels



(c) Loss robustification uses a new, robust loss function to learn the decision boundary.

Noise Matrix Estimation

- Learn a transition matrix from entire noisy dataset
- Transition Matrix could flip incorrect prediction results made by the model trained on noisy training datasets.

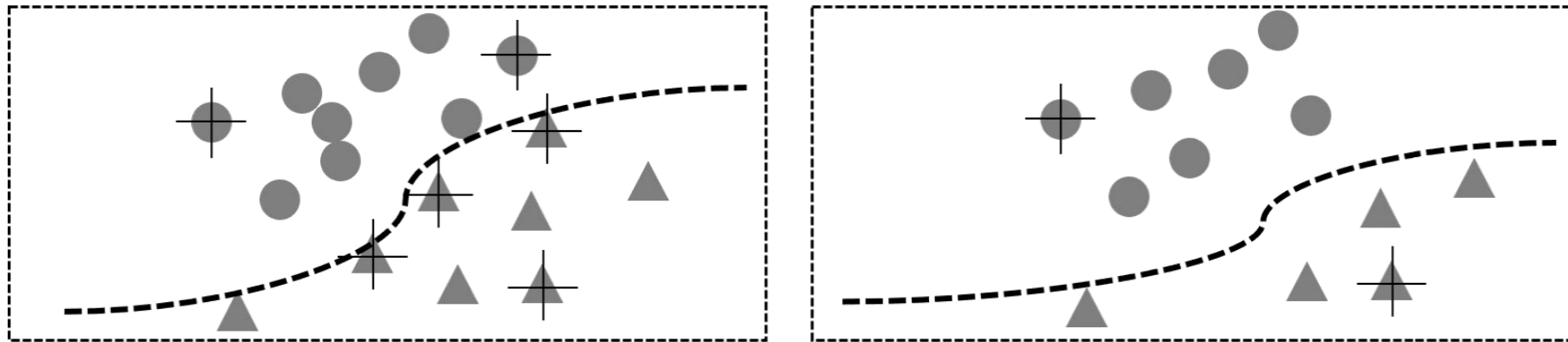


(d) Noise matrix estimation learns a transition matrix (“**T**”) and then uses it to correct the inaccurate decision boundary.

Sample Selection

- Minimize the impact of noisy labels
 - Identify incorrectly labeled samples
 - Eliminate or downplay the noisy samples from model training procedure

● ▲ Correct



(a) Sample selection selects samples possibly correctly labeled and then learns the decision boundary.

Evaluation of Existing Methods

Dataset

- Windows PE dataset
- Android Dataset

- How do we know noise rate?

TABLE 1: Statistics of the Windows PE malware dataset.

ID	Family	# of Samples	# of Training samples	Noise rate in training set
0	Benign	500	400	0.000
1	VirLock	900	800	0.005
2	WannaCry	920	820	0.002
3	Upatre	440	340	0.032
4	Cerber	1044	944	0.156
5	Urelas	572	472	0.011
6	WinActivator	166	66	0.106
7	Pykspa	744	644	0.019
8	Ramnit	324	224	0.295
9	Gamarue	608	508	0.750
10	InstallMonster	299	199	0.472
11	Locky	157	57	0.544

TABLE 2: The statistics of the Android malware dataset.

ID	Family	# of Samples	# of Training samples	Noise rate in training set
0	Smserg	2687	2587	0.040
1	Benign	4683	4583	0.699
2	Autoins	200	100	0.150
3	Jiagu	678	578	0.235
4	Shedun	10867	10767	0.548
5	Wapron	857	757	0.136
6	Dnotua	136	36	0.583
7	Hiddad	185	85	0.000
8	Secneo	203	103	0.223
9	Triada	299	199	0.005
10	Secapk	155	55	0.055
11	Smspay	281	181	0.028
12	Qappusin	158	58	0.000

Windows PE Dataset

- We assume the labels are 100% accurate (500 benign, 6174 malicious)
 - Obtained from security lab
 - Carefully analyzed by at least three security analyst with 5+ years of experience
- Generation of noisy labels
 - Used VirusTotal
 - Upload all executables to VirusTotal (discovered 11.38% of PE files provided had at least one wrong label by a vendor)
 - Randomly change the label to what the vendor provided (noisy label)

TABLE 1: Statistics of the Windows PE malware dataset.

ID	Family	# of Samples	# of Training samples	Noise rate in training set
0	Benign	500	400	0.000
1	VirLock	900	800	0.005
2	WannaCry	920	820	0.002
3	Upatre	440	340	0.032
4	Cerber	1044	944	0.156
5	Urelas	572	472	0.011
6	WinActivator	166	66	0.106
7	Pykspa	744	644	0.019
8	Ramnit	324	224	0.295
9	Gamarue	608	508	0.750
10	InstallMonster	299	199	0.472
11	Locky	157	57	0.544

Android Dataset

- VirusShare2018, 4683 benign 16706 android malware
- Difficult to check label correctness in this case
- Generation of noisy labels
 - Automatically obtained noisy labels through Virus Total
 - Upload to VT, correct label: majority vote between vendors
 - Noisy Label: from one vendor (Ikarus)

TABLE 2: The statistics of the Android malware dataset.

ID	Family	# of Samples	# of Training samples	Noise rate in training set
0	Smserg	2687	2587	0.040
1	Benign	4683	4583	0.699
2	Autoins	200	100	0.150
3	Jiagu	678	578	0.235
4	Shedun	10867	10767	0.548
5	Wapron	857	757	0.136
6	Dnotua	136	36	0.583
7	Hiddad	185	85	0.000
8	Secneo	203	103	0.223
9	Triada	299	199	0.005
10	Secapk	155	55	0.055
11	Smspays	281	181	0.028
12	Qappusin	158	58	0.000

Models to Evaluate

- SOTA, Representative Method (via citation)

TABLE 3: Summary of our selected noise learning methods.

Category	Representative Method	State-of-the-art Method
Sample selection	Coteaching+[12]	Mentormix[6]
Label sanitization	Bootstrap[27]	LRT[28]
Loss robustification	GCE[8]	ELR[9]
Noise matrix estimation	Noise-adaption[10]	LIO[13]

Windows PE Evaluation Results

Methods	Average (%)			Class-11 Locky (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	93.08/0.25	93.65/0.51	92.57/0.36	44.33/3.78	96.34/4.22	60.48/3.22
Coteaching+	87.39/0.65 $p = 0.999$	92.24/0.39 $p = 0.998$	87.90/1.13 $p = 0.999$	2.00/2.89 $p = 0.999$	99.44/1.84 $p = 0.055$	3.92/3.16 $p = 0.996$
Mentormix	92.34/0.10 $p = 0.997$	93.40/0.14 $p = 0.812$	92.32/0.28 $p = 0.322$	41.17/4.13 $p = 0.767$	95.24/1.34 $p = 0.505$	57.82/2.14 $p = 0.434$
Bootstrap	92.92/0.30 $p = 0.876$	93.33/0.23 $p = 0.829$	92.27/1.91 $p = 0.503$	46.33/5.15 $p = 0.260$	92.71/2.69 $p = 0.858$	61.03/2.49 $p = 0.191$
LRT	92.52/0.21 $p = 0.995$	93.38/0.24 $p = 0.815$	92.15/0.18 $p = 0.817$	42.50/1.50 $p = 0.887$	93.85/4.67 $p = 0.753$	59.34/2.47 $p = 0.557$
Noise-adaption	92.65/0.24 $p = 0.978$	93.25/0.30 $p = 0.861$	92.18/0.22 $p = 0.488$	40.83/1.46 $p = 0.968$	90.04/3.68 $p = 0.968$	56.18/2.47 $p = 0.687$
LIO	92.30/0.35 $p = 0.990$	93.21/0.25 $p = 0.879$	92.01/0.15 $p = 0.435$	38.00/4.04 $p = 0.910$	90.12/4.23 $p = 0.958$	53.42/1.45 $p = 0.956$
GCE	92.15/0.27 $p = 0.998$	93.40/0.32 $p = 0.809$	91.74/0.65 $p = 0.962$	36.33/7.40 $p = 0.926$	94.85/3.11 $p = 0.705$	52.38/5.95 $p = 0.991$
ELR	91.84/0.18 $p = 0.999$	93.27/0.24 $p = 0.882$	91.20/0.27 $p = 0.998$	34.50/2.75 $p = 0.999$	96.39/3.37 $p = 0.493$	50.01/2.70 $p = 0.997$

Android Evaluation Results

Methods	Average (%)			Class-6 Dnotua (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	73.00/0.35	78.65/2.21	69.96/0.54	2.67/5.53	19.79/34.31	4.63/9.52
Coteaching+	72.73/0.61 $p = 0.895$	75.83/1.76 $p = 0.949$	69.44/0.41 $p = 0.963$	0.00/0.00 $p = 0.835$	0.00/0.00 $p = 0.873$	0.00/0.00 $p = 0.837$
Mentormix	75.52/0.24 $p = 0.001$	78.20/3.12 $p = 0.645$	72.15/0.88 $p = 0.003$	3.75/6.49 $p = 0.767$	20.24/31.24 $p = 0.563$	6.30/10.18 $p = 0.620$
Bootstrap	72.24/0.42 $p = 0.993$	77.41/3.94 $p = 0.738$	68.97/0.70 $p = 0.999$	0.33/0.47 $p = 0.817$	16.67/23.57 $p = 0.591$	0.65/0.92 $p = 0.816$
LRT	72.07/0.49 $p = 0.991$	79.02/3.15 $p = 0.416$	69.15/0.87 $p = 0.933$	2.50/5.59 $p = 0.517$	15.62/34.94 $p = 0.574$	4.31/9.64 $p = 0.519$
Noise-adaption	74.73/0.80 $p = 0.003$	78.33/3.01 $p = 0.605$	71.62/0.82 $p = 0.004$	5.50/7.46 $p = 0.132$	39.17/41.87 $p = 0.114$	9.40/12.63 $p = 0.138$
LIO	72.56/0.41 $p = 0.838$	80.02/2.83 $p = 0.233$	72.78/1.00 $p = 0.001$	3.00/5.86 $p = 0.468$	25.00/38.19 $p = 0.417$	5.24/10.09 $p = 0.466$
GCE	72.67/0.95 $p = 0.710$	78.90/3.98 $p = 0.459$	69.90/0.87 $p = 0.539$	5.00/7.00 $p = 0.314$	32.29/45.70 $p = 0.333$	8.66/12.24 $p = 0.314$
ELR	70.99/0.12 $p = 0.999$	76.46/0.77 $p = 0.967$	67.62/0.25 $p = 0.999$	0.00/0.00 $p = 0.835$	0.00/0.00 $p = 0.873$	0.00/0.00 $p = 0.837$

Hypothesis

- #1 Performance degradation is due to highly skewed dataset.
(Locky Dnotua, both are smallest class in the dataset)
- #2 High noise rate reduces the number of clean samples useful for classifier training
(Locky Dnotua have high noise rate 54%, 58.3%)

Hypothesis Test

- Synthetic Dataset form BODMAS (57293 samples 581 families)
- #1 Skewed
 - Remove proportion of families obtain two dataset with different imbalance (20x and 100x)
- #2 High Noise Rate
 - Randomly change the label for each dataset (30% and 60%)

Methods	Noise rate 0.3 and imbalance ratio 20x						Noise rate 0.3 and imbalance ratio 100x					
	Overall cls (%)			Rare cls (%)			Overall cls (%)			Rare cls (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	75.55/1.38	82.94/3.78	73.46/0.90	57.89/3.69	85.74/6.65	63.03/1.20	70.83/0.87	76.29/0.85	67.80/1.19	48.45/1.44	79.47/0.62	56.77/0.24
Coteaching+	74.49/2.56 $p = 0.738$	78.23/1.45 $p = 0.951$	72.79/1.25 $p = 0.683$	56.71/3.08 $p = 0.715$	77.92/2.50 $p = 0.952$	62.09/1.03 $p = 0.871$	69.25/0.92 $p = 0.977$	74.63/3.38 $p = 0.863$	67.02/1.95 $p = 0.817$	45.57/3.68 $p = 0.891$	75.96/7.39 $p = 0.834$	53.85/4.53 $p = 0.898$
MentorMix	77.58/0.70 $p = 0.043$	80.80/1.05 $p = 0.715$	74.15/1.13 $p = 0.052$	62.77/1.14 $p = 0.058$	79.52/1.02 $p = 0.893$	63.85/1.02 $p = 0.082$	72.58/0.56 $p = 0.020$	75.25/0.24 $p = 0.998$	67.51/0.10 $p = 0.738$	48.54/1.06 $p = 0.430$	78.86/0.83 $p = 0.975$	56.73/1.02 $p = 0.481$
Bootstrap	75.17/1.43 $p = 0.624$	78.80/1.05 $p = 0.934$	72.78/1.13 $p = 0.826$	58.58/3.00 $p = 0.398$	77.55/1.79 $p = 0.940$	62.11/2.10 $p = 0.992$	70.30/0.66 $p = 0.970$	76.45/0.74 $p = 0.362$	67.55/1.24 $p = 0.637$	49.24/1.18 $p = 0.947$	79.65/0.56 $p = 0.082$	56.51/1.46 $p = 0.662$
LRT	73.17/2.16 $p = 0.927$	78.41/0.61 $p = 0.951$	72.30/1.11 $p = 0.928$	58.30/3.66 $p = 0.398$	78.25/1.10 $p = 0.934$	62.33/0.88 $p = 0.776$	65.17/5.24 $p = 0.992$	74.18/3.94 $p = 0.871$	64.42/4.93 $p = 0.923$	40.01/8.61 $p = 0.994$	77.63/4.04 $p = 0.809$	51.71/6.40 $p = 0.935$
Noise-adaption	77.75/0.66 $p = 0.016$	80.44/1.06 $p = 0.846$	74.13/1.33 $p = 0.066$	62.50/0.98 $p = 0.028$	79.02/2.68 $p = 0.943$	63.73/1.56 $p = 0.096$	72.56/1.41 $p = 0.006$	75.59/0.81 $p = 0.846$	67.30/1.51 $p = 0.691$	50.81/2.90 $p = 0.075$	79.21/0.57 $p = 0.975$	56.76/0.79 $p = 0.504$
LIO	75.72/1.78 $p = 0.825$	81.19/2.74 $p = 0.642$	72.73/0.75 $p = 0.883$	60.19/1.04 $p = 0.114$	84.40/6.52 $p = 0.479$	63.61/2.44 $p = 0.238$	70.06/2.34 $p = 1.000$	75.12/2.35 $p = 0.852$	65.80/1.12 $p = 0.992$	46.25/5.68 $p = 0.798$	79.65/0.30 $p = 0.306$	52.72/0.28 $p = 0.987$
GCE	76.98/1.49 $p = 0.113$	81.52/3.26 $p = 0.594$	73.68/1.74 $p = 0.287$	59.08/4.09 $p = 0.350$	82.89/5.96 $p = 0.614$	62.67/2.38 $p = 0.540$	70.83/2.21 $p = 0.504$	73.25/7.41 $p = 0.793$	66.94/3.21 $p = 0.704$	47.58/3.32 $p = 0.728$	72.31/14.53 $p = 0.836$	53.49/6.24 $p = 0.848$
ELR	74.28/0.69 $p = 0.953$	84.60/1.20 $p = 0.077$	72.05/1.38 $p = 0.899$	53.58/2.10 $p = 0.982$	92.02/1.57 $p = 0.027$	60.82/1.97 $p = 0.944$	66.84/1.66 $p = 0.998$	65.70/0.13 $p = 1.000$	61.00/0.40 $p = 0.999$	42.33/3.37 $p = 0.987$	59.46/0.57 $p = 1.000$	44.90/0.77 $p = 1.000$
Methods	Noise rate 0.6 and imbalance ratio 20x						Noise rate 0.6 and imbalance ratio 100x					
	Overall cls (%)			Rare cls (%)			Overall cls (%)			Rare cls (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	68.04/1.77	70.38/6.31	63.68/4.12	48.93/4.03	73.16/9.31	53.39/7.05	65.57/0.87	71.67/3.30	61.11/0.89	44.39/1.94	76.67/7.45	51.02/1.56
Coteaching+	61.20/4.81 $p = 0.993$	58.92/5.03 $p = 0.997$	54.73/4.60 $p = 0.998$	33.64/6.49 $p = 0.998$	55.50/9.97 $p = 0.997$	39.84/8.42 $p = 0.997$	50.85/7.03 $p = 0.997$	37.47/4.82 $p = 1.000$	41.93/4.07 $p = 0.999$	10.37/14.67 $p = 0.998$	16.38/7.34 $p = 1.000$	15.33/6.87 $p = 1.000$
MentorMix	69.11/0.74 $p = 0.129$	75.05/1.05 $p = 0.934$	63.28/1.03 $p = 0.826$	51.40/2.28 $p = 0.376$	77.55/3.79 $p = 0.151$	58.11/2.10 $p = 0.124$	55.28/1.50 $p = 0.998$	76.45/0.74 $p = 0.362$	57.55/1.24 $p = 0.999$	19.96/2.34 $p = 0.999$	77.58/4.01 $p = 0.082$	30.51/3.46 $p = 0.999$
Bootstrap	68.94/1.40 $p = 0.174$	71.12/4.69 $p = 0.431$	64.81/2.92 $p = 0.328$	50.03/5.04 $p = 0.370$	71.35/8.21 $p = 0.602$	55.56/7.42 $p = 0.346$	65.83/1.50 $p = 0.316$	72.73/0.67 $p = 0.253$	60.89/1.10 $p = 0.669$	45.66/2.09 $p = 0.182$	79.95/0.07 $p = 0.184$	50.30/1.21 $p = 0.784$
LRT	64.97/2.70 $p = 0.982$	70.83/2.85 $p = 0.453$	61.18/4.68 $p = 0.759$	48.40/3.06 $p = 0.643$	78.15/1.91 $p = 0.168$	54.61/2.64 $p = 0.380$	55.51/5.83 $p = 0.992$	57.23/6.02 $p = 0.993$	47.04/5.76 $p = 0.998$	26.78/11.11 $p = 0.994$	59.15/10.55 $p = 0.962$	29.60/10.03 $p = 0.997$
Noise-adaption	70.62/3.54 $p = 0.130$	77.15/2.84 $p = 0.021$	66.75/2.11 $p = 0.141$	53.65/7.79 $p = 0.155$	83.70/4.08 $p = 0.010$	58.07/3.34 $p = 0.146$	66.20/1.79 $p = 0.196$	72.04/3.48 $p = 0.443$	62.01/1.18 $p = 0.120$	44.82/2.55 $p = 0.380$	75.80/7.18 $p = 0.563$	53.60/2.67 $p = 0.101$
LIO	68.54/2.97 $p = 0.390$	73.43/2.67 $p = 0.167$	64.63/2.40 $p = 0.311$	48.73/6.64 $p = 0.517$	77.28/3.51 $p = 0.171$	57.39/3.04 $p = 0.133$	52.54/5.26 $p = 1.000$	56.90/7.90 $p = 0.997$	45.92/6.74 $p = 0.998$	24.53/9.29 $p = 0.997$	57.12/15.47 $p = 0.980$	27.23/14.41 $p = 0.993$
GCE	57.38/4.32 $p = 0.996$	53.55/6.63 $p = 0.993$	50.76/6.05 $p = 0.992$	29.01/7.37 $p = 0.997$	45.59/10.18 $p = 0.994$	32.45/7.67 $p = 0.995$	47.15/0.35 $p = 1.000$	27.53/1.75 $p = 1.000$	33.92/1.02 $p = 1.000$	0/0 $p = 1.000$	0/0 $p = 1.000$	0/0 $p = 1.000$
ELR	68.19/1.93 $p = 0.430$	68.62/5.42 $p = 0.727$	62.58/3.58 $p = 0.719$	43.55/4.02 $p = 0.985$	66.42/9.61 $p = 0.914$	48.51/5.95 $p = 0.917$	63.65/0.53 $p = 0.994$	60.18/5.11 $p = 0.998$	58.69/1.41 $p = 0.999$	36.13/1.07 $p = 1.000$	50.00/10.00 $p = 0.999$	40.00/2.46 $p = 1.000$

Malware != Images

- Highly Skewed and High Noise Rate in Malware dataset require different noise learning approach
- MORSE
- Semi Supervised learning method + Sample re-weighting Mechanism

Semi Supervised Learning: FixMatch

- Tailored for Image Recognition
- Partition labeled and unlabeled datasets
- Weakly augment (flip, shift) unlabeled data
- Predict labels for unlabeled data (pseudo labels)
- Calculate Both supervised, unsupervised loss
- Calculate Final Loss and update model's weights

Algorithm 1: FixMatch – the state-of-the-art semi-supervised learning algorithm.

```
1 Input: labeled dataset  $\mathcal{X}$ , unlabeled dataset  $\mathcal{U}$ ,  
   number of total training epochs  $K$ , confidence  
   threshold  $\tau$ , unsupervised loss weight  $\lambda$ .  
2 Initialization: Initialize the weights  $\Theta$  for the  
   model  $f(\cdot)$  randomly.  
3 for  $k = 0, 1, 2, \dots, K$  do  
4   for  $iter = 1, 2, \dots, num\_batches$  do  
5     From  $\mathcal{X}$ , draw a mini-batch  
        $\{(\mathbf{x}_b, y_b) : b \in (1, \dots, B)\}$   
6     From  $\mathcal{U}$ , draw a mini-batch  
        $\{\mathbf{u}_b : b \in (1, \dots, B_\mu)\}$   
7     for  $b = 1, 2, \dots, B_\mu$  do  
8        $\mathbf{q}_b = f(g(\mathbf{u}_b); \Theta)$   
9        $\hat{q}_b = \arg \max(\mathbf{q}_b)$   
10    end  
11     $\mathcal{L}_s = \frac{1}{B} \sum_{b=1}^B H(y_b, f(g(\mathbf{x}_b)))$   
12     $\mathcal{L}_u = \frac{1}{B_\mu} \sum_{b=1}^{B_\mu} \mathbf{I}(\max(\mathbf{q}_b) >$   
        $\tau) H(\hat{q}_b, f(h(\mathbf{u}_b)))$   
13     $\mathcal{L} = \mathcal{L}_s + \lambda \mathcal{L}_u$   
14    Update the model's weights  $\Theta$  by  
       minimizing the loss function  $\mathcal{L}$   
15  end  
16 end  
17 Output: the well trained model  $f(\cdot; \Theta)$ .
```

MORSE

- FixMatch -> MORSE
- Change design to suit Malware Classification (Augmentation)
- Add Sample Re-weighting

Algorithm 2: Proposed learning algorithm. The customized and extended parts are highlighted.

- 1 **Input:** imbalanced noisy training dataset D ,
number of total training epochs K , labeled data's
proportion d , starting re-weighting epoch T_d ,
confidence threshold τ , unsupervised loss weight
 λ , learning rate α .
- 2 **Initialization:** Initialize the weights Θ for the
model $f(\cdot)$ by using entire dataset D with only a
few epochs.
- 3 **for** $k = 0, 1, 2, \dots, K$ **do**
- 4 Partitioning the training dataset D into labeled
dataset \mathcal{X} and unlabeled dataset \mathcal{U} : select the
top $d\%$ examples with the least loss values
from each given class and treat them as labeled
data and the rest as unlabeled data.
- 5 **for** $iter = 1, 2, \dots, num_batches$ **do**
- 6 From \mathcal{X} , draw a mini-batch
 $\{(\mathbf{x}_b, y_b) : b \in (1, \dots, B)\}$
- 7 From \mathcal{U} , draw a mini-batch
 $\{\mathbf{u}_b : b \in (1, \dots, B_\mu)\}$
- 8 **for** $b = 1, 2, \dots, B_\mu$ **do**
- 9 // Use Eqn. (3) to perform weak
 augmentation against \mathbf{u}_b
- 10 $\mathbf{q}_b = f(g(\mathbf{u}_b); \Theta)$
- 11 $\hat{q}_b = \arg \max(\mathbf{q}_b)$
- 12 **end**

- 13 **if** $k < T_d$ **then**
- 14 // Use Eqn. (3) to perform weak
 augmentation against \mathbf{x}_b
- 15 $\mathcal{L}_s = \frac{1}{B} \sum_{b=1}^B H(y_b, f(g(\mathbf{x}_b)))$
- 16 // Use Eqn. (3) to perform strong
 augmentation against \mathbf{u}_b
- 17 $\mathcal{L}_u = \frac{1}{B_\mu} \sum_{b=1}^{B_\mu} \mathbf{I}(\max(\mathbf{q}_b) >$
 $\tau) H(\hat{q}_b, f(h(\mathbf{u}_b)))$
- 18 **else**
- 19 // Calculate the weight of each training
 sample using Eqn. (4)
- 20 // Use Eqn. (3) to perform weak
 augmentation against \mathbf{x}_b
- 21 $\mathcal{L}_s = \frac{1}{B} \sum_{b=1}^B w_b H(y_b, f(g(\mathbf{x}_b)))$
- 22 // Use Eqn. (3) to perform strong
 augmentation against \mathbf{u}_b
- 23 $\mathcal{L}_{u1} = \frac{1}{B_\mu} \sum_{b=1}^{B_\mu} \mathbf{I}(\max(\mathbf{q}_b) > \tau)$
- 24 $\mathcal{L}_u = \mathcal{L}_{u1} + w_b H(\hat{q}_b, f(h(\mathbf{u}_b)))$
- 25 **end**
- 26 $\mathcal{L} = \mathcal{L}_s + \lambda \mathcal{L}_u$
- 27 Update the model's weights Θ by
 minimizing the loss function \mathcal{L}
- 28 Optional: decay the learning rate α
- 29 **end**
- 30 **end**
- 31 **Output:** the well trained model $f(\cdot; \Theta)$.

Selecting a Label

- Assume a pretrained model (At initialization)
- Each Epoch we determine labeled and unlabeled data
 - Calculate loss, top $d\%$ samples with least loss -> labeled dataset

Weak & Strong augmentation

- Mask vector $[m_1, \dots, m_d]^T \in \mathbb{R}^d$ Bernoulli distribution
- Replace features with from $\bar{\mathbf{x}}$
- Where, $\bar{\mathbf{x}}$ is sampled from the entire dataset

$$\tilde{\mathbf{x}} = \mathbf{x} \odot \mathbf{m} + (1 - \mathbf{m}) \odot \bar{\mathbf{x}}.$$

Sample Re-Weighting

- Handle Class imbalance (weigh the loss differently depending on whether they belong to majority or minority class)
- Samples associated with minor class is assigned higher weight

$$E_{n_{y_b}} = \frac{1 - \beta^{n_{y_b}}}{1 - \beta}. \quad (4)$$

n_{y_b} is the number of samples in class y_b , and $\beta \in [0, 1)$ is a hyperparameter. Using its inverse (*i.e.*, $w_b = \frac{1}{E_{n_{y_b}}}$), we can extend the loss shown in Equation (1) and (2) as

$$\begin{aligned} \mathcal{L}_s &= \frac{1}{B} \sum_{b=1}^B w_b H(y_b, f(g(\mathbf{x}_b))); \\ \mathcal{L}_u &= \frac{1}{B_\mu} \sum_{b=1}^{B_\mu} \mathbf{I}(\max(\mathbf{q}_b) > \tau) w_b H(\hat{q}_b, f(h(\mathbf{u}_b))). \end{aligned} \quad (5)$$

Synthetic Dataset

TABLE 7: The testing performance of the models learned from MORSE and Vanilla method on the synthetic dataset.

Methods	Noise rate 0.3 and imbalance ratio 20x						Noise rate 0.3 and imbalance ratio 100x					
	Overall cls (%)			Rare cls (%)			Overall cls (%)			Rare cls (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	75.55/1.38	82.94/3.78	73.46/0.90	57.89/3.69	85.74/6.65	63.03/1.20	70.83/0.87	76.29/0.85	67.80/1.19	48.45/1.44	79.47/0.62	56.77/0.24
Vanilla DNN w re-weighting	78.87/0.30 $p = 0.000$	82.64/0.88 $p = 0.335$	77.71/0.44 $p = 0.000$	68.59/2.29 $p = 0.000$	80.75/2.44 $p = 0.878$	69.63/1.03 $p = 0.000$	72.85/1.79 $p = 0.083$	75.97/0.98 $p = 0.656$	69.86/1.74 $p = 0.049$	56.81/4.63 $p = 0.009$	76.44/1.89 $p = 0.988$	60.88/2.64 $p = 0.008$
Our method w/o re-weighting	79.64/0.83 $p = 0.000$	87.81/1.68 $p = 0.011$	77.61/1.01 $p = 0.000$	64.85/2.62 $p = 0.001$	94.46/2.31 $p = 0.009$	70.00/2.70 $p = 0.000$	71.73/2.24 $p = 0.138$	73.38/4.94 $p = 0.900$	65.87/3.16 $p = 0.891$	46.84/3.98 $p = 0.752$	72.99/9.19 $p = 0.917$	50.54/5.68 $p = 0.972$
Our method	79.73/1.85 $p = 0.008$	87.11/1.82 $p = 0.039$	79.11/1.30 $p = 0.000$	67.47/3.17 $p = 0.006$	88.45/4.63 $p = 0.182$	70.76/3.42 $p = 0.002$	73.47/0.91 $p = 0.001$	78.97/2.11 $p = 0.027$	70.42/1.02 $p = 0.007$	54.40/7.97 $p = 0.092$	82.59/3.85 $p = 0.071$	62.30/5.55 $p = 0.040$
Methods	Noise rate 0.6 and imbalance ratio 20x						Noise rate 0.6 and imbalance ratio 100x					
	Overall cls (%)			Rare cls (%)			Overall cls (%)			Rare cls (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	68.04/1.77	70.38/6.31	63.68/4.12	48.93/4.03	73.16/9.31	53.39/7.05	65.57/0.87	71.67/3.30	61.11/0.89	44.39/1.94	76.67/7.45	51.02/1.56
Vanilla DNN w re-weighting	68.01/1.52 $p = 0.464$	70.65/2.95 $p = 0.448$	65.52/2.68 $p = 0.184$	58.50/2.59 $p = 0.011$	67.70/5.18 $p = 0.960$	57.18/3.31 $p = 0.106$	66.84/2.36 $p = 0.137$	67.22/2.28 $p = 0.979$	64.34/2.33 $p = 0.021$	55.85/4.36 $p = 0.001$	61.15/4.13 $p = 0.996$	56.36/4.29 $p = 0.037$
Our method w/o re-weighting	72.82/2.86 $p = 0.012$	77.24/1.69 $p = 0.050$	69.42/2.61 $p = 0.038$	54.09/3.96 $p = 0.051$	81.00/1.34 $p = 0.061$	60.26/2.85 $p = 0.041$	64.85/4.83 $p = 0.624$	65.85/5.03 $p = 0.927$	58.93/2.81 $p = 0.988$	38.65/9.12 $p = 0.895$	66.67/8.43 $p = 0.898$	42.38/6.12 $p = 0.987$
Our method	76.20/0.59 $p = 0.000$	79.45/1.69 $p = 0.016$	72.90/1.17 $p = 0.003$	61.21/1.00 $p = 0.000$	80.45/0.99 $p = 0.066$	64.25/0.20 $p = 0.009$	72.05/1.94 $p = 0.000$	74.97/3.41 $p = 0.005$	68.03/1.01 $p = 0.000$	53.82/4.98 $p = 0.005$	76.60/7.48 $p = 0.891$	58.00/3.80 $p = 0.003$

PE and Android Dataset

TABLE 8: MORSE vs. Vanilla method on the PE dataset.

Methods	Average (%)			Class-11 (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	93.08/0.25	93.65/0.51	92.57/0.36	44.33/3.78	96.34/4.22	60.48/3.22
Vanilla DNN w re-weighting	93.82/0.29 $p = 0.003$	94.01/0.30 $p = 0.102$	93.50/0.34 $p = 0.001$	55.67/2.98 $p = 0.001$	89.71/2.30 $p = 0.974$	68.61/1.70 $p = 0.001$
Our method w/o re-weighting	92.50/0.18 $p = 0.980$	93.35/0.28 $p = 0.924$	91.33/0.58 $p = 0.992$	36.50/2.75 $p = 0.995$	99.51/1.10 $p = 0.089$	51.07/5.05 $p = 0.979$
Our method	94.15/0.43 $p = 0.003$	94.14/0.57 $p = 0.061$	93.91/0.67 $p = 0.000$	65.33/4.20 $p = 0.000$	90.18/2.08 $p = 0.984$	76.74/1.22 $p = 0.000$

TABLE 9: MORSE vs. Vanilla on the Android dataset.

Methods	Average (%)			Class-6 (%)		
	Accuracy	Precision	F_1	Accuracy	Precision	F_1
Vanilla DNN	73.00/0.35	78.65/2.21	69.96/0.54	2.67/5.53	19.79/34.31	4.63/9.52
Vanilla DNN w re-weighting	77.32/1.72 $p = 0.002$	82.20/3.73 $p = 0.049$	75.68/2.58 $p = 0.003$	20.83/7.99 $p = 0.007$	80.39/9.40 $p = 0.003$	32.03/8.99 $p = 0.005$
Our method w/o re-weighting	74.58/0.86 $p = 0.004$	77.33/4.47 $p = 0.691$	71.00/0.96 $p = 0.015$	5.17/7.31 $p = 0.309$	32.35/45.79 $p = 0.347$	8.91/12.60 $p = 0.311$
Our method	79.76/1.02 $p = 0.000$	80.37/1.99 $p = 0.088$	78.72/1.22 $p = 0.000$	27.67/8.01 $p = 0.000$	58.54/18.90 $p = 0.039$	35.80/7.88 $p = 0.000$