

Weight-based Semantic Testing Approach for Neural Networks

Amany Alshareef, Nicolas Berthier, Sven Schewe, Xiaowei Huang

Department of Computer Science
University of Liverpool, UK



UNIVERSITY OF
LIVERPOOL

Testing goal

- Evidence that the system meets its requirements
- Evidence that the system is error-free

Machine learned software !!

- Data-driven NOT requirements-driven
- Do NOT have a specific control-flow structure
- Most testing techniques propose structural coverage
- Tend at transforming in the input data



Real world high-dimensional data lie on low-dimensional manifolds embedded within the high-dimensional.

→ Not thinking about input ~~domain~~, instead where the data lie and model that



The research presents a testing approach for neural networks that leverages the learned representations and feature importance to evaluate the test data coverage.



Analyse the latent features
learnt by the model



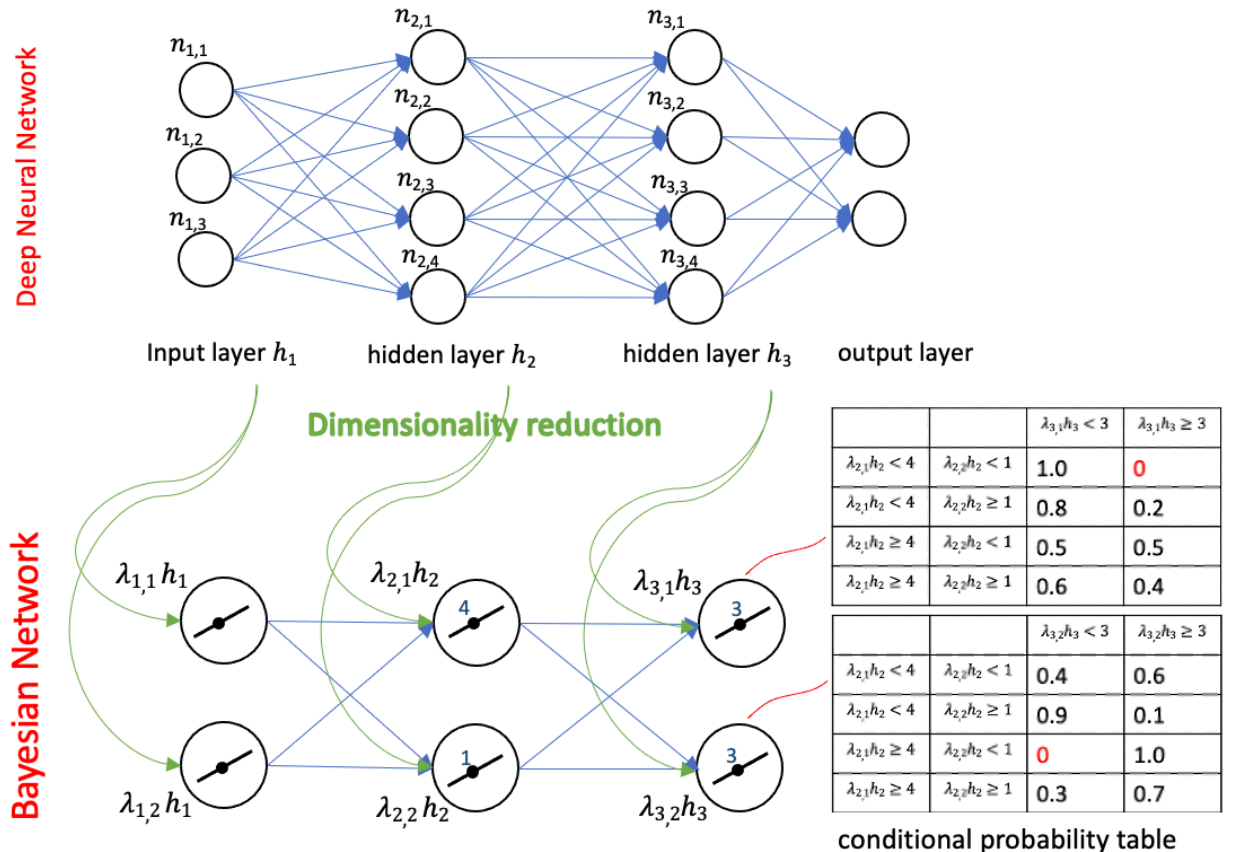
Generate additional
test cases

Bayesian Network Abstraction Model

A dimensionality reduction technique using feature extraction algorithms to abstract the behaviour of a neural network into a BN.

Constructing a Bayesian Network:

- I. **Hidden features extraction**
Map from a high-dimensional space into a feature space.
- II. **Feature space discretization**
Discretise each feature component into finite feature intervals.
- III. **Probability tables construction**
 - Represent the probabilistic distribution of each extracted feature with a node in the BN.
 - Associate each feature with a marginal or a conditional probability table.



1 BN-based Latent Feature Analysis

Estimate the **importance** of a neural network's **latent features** by analysing an associated **Bayesian network's sensitivity** to **distributional shifts**

1. Probability calculation for input sample under BN

Perform the feature projection and discretisation step to input sample to obtain the associated feature intervals, and then calculate their probability belonging to the BN distribution.

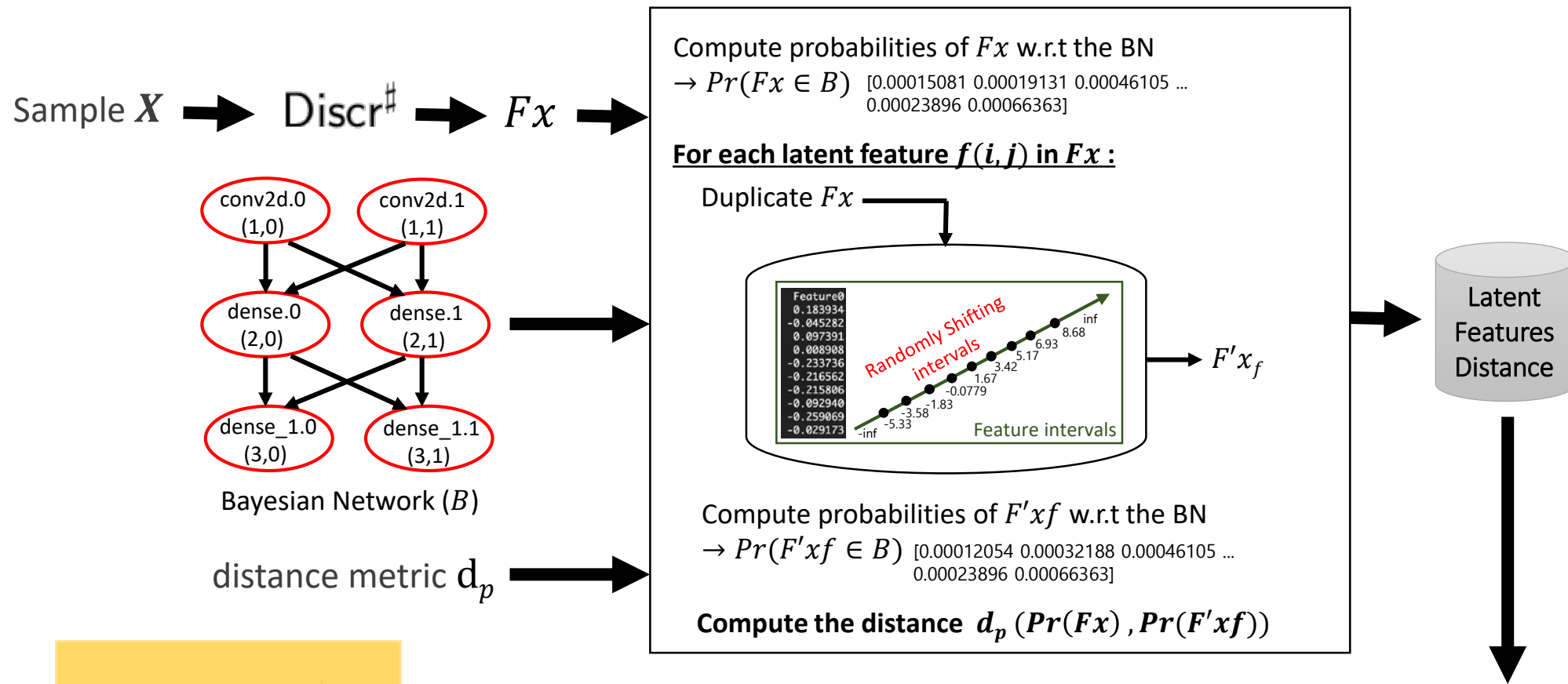
2. Latent features perturbation

For each latent feature, randomly shifting its intervals in a selected feature space.

3. Distance computation

Compute the distance between the original probability vector and the probability vector obtained from the perturbed features.

The proposed BN analysis technique to compute the sensitivity of extracted latent features



$$W_{i,j} = \frac{\delta(P_{ref}, P'_f)}{\sum_{f \in F\#} \delta(P_{ref}, P'_f)}$$

$W_{(1,1)} = 0.192$
 $W_{(2,1)} = 0.172$

Higher distribution change → Higher importance score

distance perturbed feature	d_{L_1}	d_{L_2}	d_{L_∞}	d_{JS}	d_{corr}	d_{cos}	d_{MSE}	d_{RMSE}	d_{MAE}	d_{AF}
(1, 0)	150	0.726	0.00956	0.224	0.142	0.114	0.000000879	0.000937	0.000249	0.278
(1, 1)	340	1.18	0.00989	0.353	0.448	0.361	0.00000232	0.00152	0.000567	0.735
(2, 0)	325	1.09	0.00946	0.365	0.332	0.267	0.00000198	0.00141	0.000541	0.625
(2, 1)	360	1.16	0.0103	0.393	0.395	0.323	0.00000224	0.00150	0.000600	0.710
(3, 0)	276	0.880	0.00889	0.258	0.170	0.137	0.00000129	0.00114	0.000460	0.408
(3, 1)	315	1.07	0.00960	0.324	0.318	0.264	0.00000192	0.00139	0.000525	0.608

2 Weight-based Testing Metric

Transforming the traditional binary coverage approach to a weighted probability problem and define our coverage metric based on the latent features importance.

- Weight-based Feature Coverage

$$\text{WFCov}(\mathcal{B}_{\mathcal{N},\mathcal{X}}) \stackrel{\text{def}}{=} \sum_{(f_{i,j}^\#) \in V_{\mathcal{N},\mathcal{X}}} w_{(f_{i,j}^\#)} \cdot \frac{|\{f_{i,j}^{\#k} \in \mathbb{F}_{i,j}^\# \mid \mathcal{P}_i(f_{i,j}^{\#k}) \geq \varepsilon\}|}{|\mathbb{F}_{i,j}^\#|}$$

- Weight-based Feature Dependence Coverage

$$\text{WFdCov}(\mathcal{B}_{\mathcal{N},\mathcal{X}}) \stackrel{\text{def}}{=} \sum_{(f_{i,j}^\#) \in V_{\mathcal{N},\mathcal{X}}^+} w_{(f_{i,j}^\#)} \cdot \frac{\left| \begin{array}{l} (f_{i,j}^{\#k}, F_{i-1}^\#) \in \mathcal{CP}_i(f_{i,j}^{\#k} | F_{i-1}^\#) \geq \varepsilon \\ \mathbb{F}_{i,j}^\# \times \mathbb{F}_{i-1}^\# \quad \vee \quad \mathcal{P}_i(f_{i,j}^{\#k}) < \varepsilon \end{array} \right|}{|\mathbb{F}_{i,j}^\# \times \mathbb{F}_{i-1}^\#|}$$

The two feature metrics are combined to produce the generalised weight feature coverage.

$$\text{WFCovTot}(\mathcal{B}_{\mathcal{N},\mathcal{X}}) = \sum_{(f_{i,j}^\#) \in V_{\mathcal{N},\mathcal{X}}} w_{(f_{i,j}^\#)} \cdot \begin{cases} \text{WFCov}_{(f_{i,j}^\#)} & \text{if } i = 1, \\ \frac{1}{2} \left(\text{WFCov}_{(f_{i,j}^\#)} + \text{WFdCov}_{(f_{i,j}^\#)} \right) & \text{otherwise.} \end{cases}$$

2 Weight-based Testing Metric

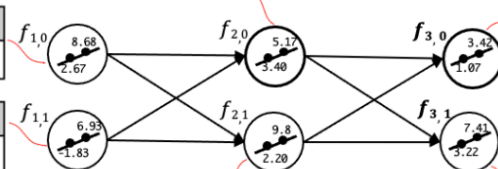
		$f_{2,0} < 3$	$3 \leq f_{2,0} \leq 5$	$f_{2,0} > 5$
$f_{1,0} < 2$	$f_{1,1} < -1$	0.44	0.24	0.32
$2 \leq f_{1,0} \leq 8$	$f_{1,1} < -1$	0.75	0.25	0.00
...
$f_{1,0} > 8$	$f_{1,1} > 6$	0.16	0.75	0.09
marginal distribution		0.343	0.207	0.450
node weight		0.1976		

		$f_{3,0} < 3$	$3 \leq f_{3,0} \leq 5$	$f_{3,0} > 5$
$f_{2,0} < 3$	$f_{2,1} < 2$	0.90	0.05	0.05
$f_{2,0} < 3$	$2 \leq f_{2,1} \leq 9$	0.10	0.70	0.20
$f_{2,0} < 3$	$f_{2,1} > 9$	0.80	0.15	0.05
$3 \leq f_{2,0} \leq 5$	$f_{2,1} < 2$	0.20	0.80	0.00
$3 \leq f_{2,0} \leq 5$	$2 \leq f_{2,1} \leq 9$	0.35	0.34	0.31
$3 \leq f_{2,0} \leq 5$	$f_{2,1} > 9$	0.40	0.30	0.30
$f_{2,0} > 5$	$f_{2,1} < 2$	0.38	0.53	0.09
$f_{2,0} > 5$	$2 \leq f_{2,1} \leq 9$	0.45	0.03	0.52
$f_{2,0} > 5$	$f_{2,1} > 9$	0.47	0.02	0.51
marginal distribution		0.453	0.323	0.224
node weight		0.1730		

		$f_{3,1} < 3$	$3 \leq f_{3,1} \leq 7$	$f_{3,1} > 7$
$f_{2,0} < 3$	$f_{2,1} < 2$	0.4	0.6	0.0
$f_{2,0} < 3$	$2 \leq f_{2,1} \leq 9$	0.1	0.2	0.7
$f_{2,0} < 3$	$f_{2,1} > 9$	0.55	0.06	0.39
$3 \leq f_{2,0} \leq 5$	$f_{2,1} < 2$	0.15	0.8	0.05
$3 \leq f_{2,0} \leq 5$	$2 \leq f_{2,1} \leq 9$	0.3	0.3	0.4
$3 \leq f_{2,0} \leq 5$	$f_{2,1} > 9$	0.2	0.6	0.2
$f_{2,0} > 5$	$f_{2,1} < 2$	0.31	0.34	0.35
$f_{2,0} > 5$	$2 \leq f_{2,1} \leq 9$	0.9	0.0	0.1
$f_{2,0} > 5$	$f_{2,1} > 9$	0.5	0.4	0.1
marginal distribution		0.380	0.366	0.254
node weight		0.2697		

$f_{1,0} < 2$	$2 \leq f_{1,0} \leq 8$	$f_{1,0} > 8$	weight
0.518	0.150	0.332	0.0318

$f_{1,1} < -1$	$-1 \leq f_{1,1} \leq 6$	$f_{1,1} > 6$	weight
0.129	0.776	0.095	0.0992



		$f_{2,1} < 2$	$2 \leq f_{2,1} \leq 9$	$f_{2,1} > 9$
$f_{1,0} < 2$	$f_{1,1} < -1$	0.60	0.10	0.30
$f_{1,0} < 2$	$-1 \leq f_{1,1} \leq 6$	0.16	0.75	0.09
...
$f_{1,0} > 8$	$f_{1,1} > 6$	0.50	5.00	0.0
marginal distribution		0.243	0.521	0.236
node weight		0.2284		

Example 1

$\Pr(f_{3,0} < 3) \approx 0.453$, $\Pr(3 \leq f_{3,0} \leq 5) \approx 0.323$, $\Pr(f_{3,0} > 5) \approx 0.224$.

$$1 \times 0.173 = 0.173.$$

$$\mathbf{WFCov(B_N, x) = 1}$$

Example 2

$$\epsilon = 0.01$$

$$CP_{f_{3,0}} = 0.1730 / 0.8687 \cdot 26 / 27 =$$

$$CP_{f_{3,1}} = \frac{0.1917}{0.2697} / 0.8687 \cdot 25 / 27 =$$

$$\mathbf{WFdCov(B_N, x) = 0.2190 +$$

$$\mathbf{0.2532 + 0.1917 + 0.2875 =$$

$$\mathbf{0.9514\%}.$$

Concolic Test Data Generation Algorithm

Algorithm 1 Test Dataset Generation

Input:

$\mathcal{N} \leftarrow$ DNN under test

$X \leftarrow$ data set

$\mathcal{B}_{\mathcal{N}, X_{train}} \leftarrow$ abstract BN

$W_f \leftarrow$ features sensitivity weights

Output: test inputs X_0 , coverage

- 1: $X_0 \leftarrow$ sampling initial seed test inputs from X_{test}
 - 2: $\mathcal{B}_{\mathcal{N}, X_0} \leftarrow$ initialising the BN prob. tables with X_0
 - 3: $Tar_invals \leftarrow$ intervals with prob $\leq \epsilon$
 - 4: **for** $i = 1$ to max iterations **do**
 - 5: $t \leftarrow Tar_invals$ with highest weight in W_f
 - 6: select a test input $s \in X_0$
 - 7: construct an LP problem based on t
 - 8: solve the optimisation objective:
 $\min \|(n_{1,1}, \dots, n_{1,|l_1|}) - (s_{1,1}, \dots, s_{1,|l_1|})\|_\infty$
 - 9: $s' = (n_{1,1}, \dots, n_{1,|l_1|})$
 - 10: **if** s' passes the oracle **then**
 - 11: $s' \leftarrow$ newly generated test input
 - 12: **if** $f_{\mathcal{N}}(s') = f_{\mathcal{N}}(s)$ **then**
 - 13: $X_0 \leftarrow X_0 \cup \{s'\}$
 - 14: update $\mathcal{B}_{\mathcal{N}, X_0}$ probabilities
 - 15: update coverage
 - 16: **else**
 - 17: $s' \leftarrow$ adversarial input
 - 18: **end if**
 - 19: **end if**
 - 20: **end for**
-

Evaluation

1. Datasets and Models

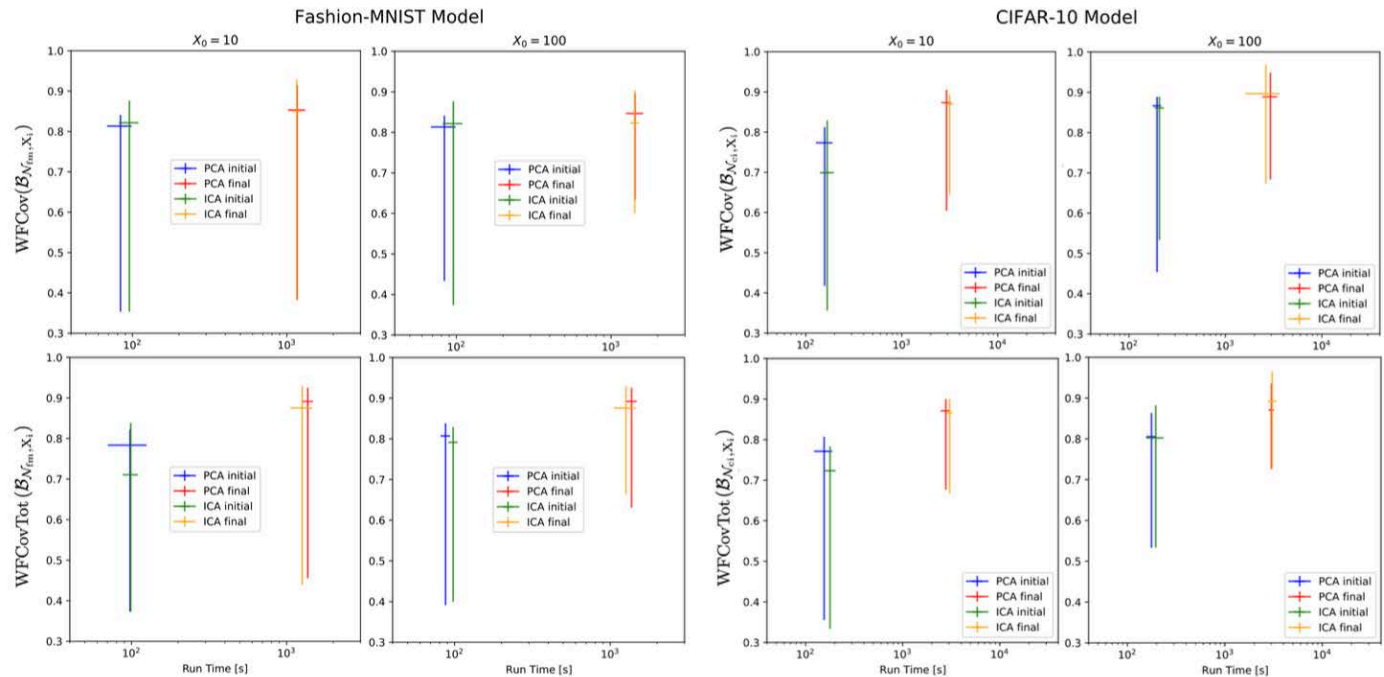
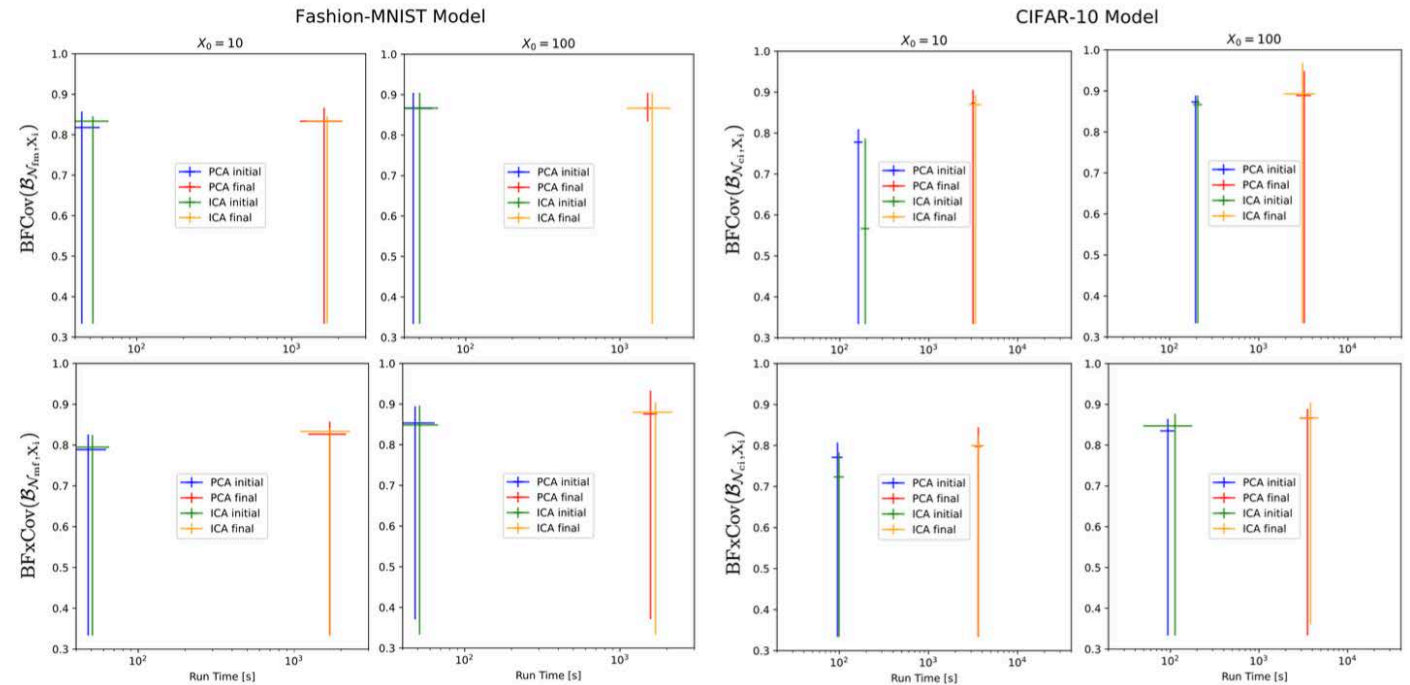
- ❑ The first model targets the **Fashion-MNIST** classification problem with 89.03% validation accuracy, and the second one targets the **CIFAR-10** dataset with 81.00% validation accuracy.
- ❑ The models are reasonably sized, with more than 10 layers, including blocks of convolutional and max-pooling layers, followed by a series of dense layers.
- ❑ The considered layers are the convolutional ReLU, 2d max pooling, and dense ReLU.

2. Experimental Setup

- ❑ Two linear feature extraction techniques were selected: PCA and ICA, with two to five numbers of extracted features for each of the abstracted layers.
- ❑ The Kernel Density Estimation (KDE) and uniform-based discretisation are considered, with varying numbers of the uniform partitions bins that are: one, three, and five.
- ❑ The extended Concolic testing tool is run on both DNNs models with a maximum 100 iterations per run. Each run is initialised with uniformly drawn test sets of 10 and 100 correctly classified inputs.



Results



The overall distribution of initial and the respective final coverage of up to 100 iterations of Concolic test case generation. X-axis indicates the run time in seconds (initial and run time). The vertical lines are the coverage, and the horizontal line on the coverage is the median.



We have introduced a weight-based semantic testing approach that measures how well the DNN is tested by focusing on the important features of the DNN using its abstracted Bayesian network.

- ✓ The developed weighted feature metrics achieved higher testing coverage than the original metrics.
- ✓ The test generation algorithm is directed to synthesise new input targeting features with higher importance scores.
- ✓ Empirically validated the applicability and effectiveness of the proposed weight metrics, which serves as a strong argument in favour of increasing the trustworthy performance of the DNN models.

THANK YOU!