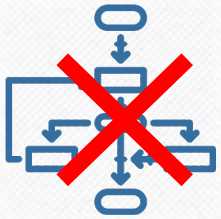# Quantifying the Importance of Latent Features in Neural Networks

**Amany Alshareef, Nicolas Berthier, Sven Schewe, Xiaowei Huang**

Department of Computer Science
University of Liverpool, UK

UNIVERSITY OF
LIVERPOOL

# Challenges in Deep Neural Network Testing

→ Deep neural networks learn by example
→ Do **not** have a specific control-flow structure

→ Most testing techniques propose structural coverage
→ Tend at transforming in the input data

The factor causing the adversarial vulnerability is the distortion in the **latent feature space**.

→ **Need to explore the internal logic of the learning models**

UNIVERSITY OF
LIVERPOOL

# Aim and Contribution

**Understand the DNN underlying decision processes**

Analyse the latent features learnt by the model

Generate additional test cases based on that

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

# Bayesian Network

A dimensionality reduction technique using feature extraction algorithms to abstracts the behaviour of a DNN.

**Constructing a Bayesian Network:**

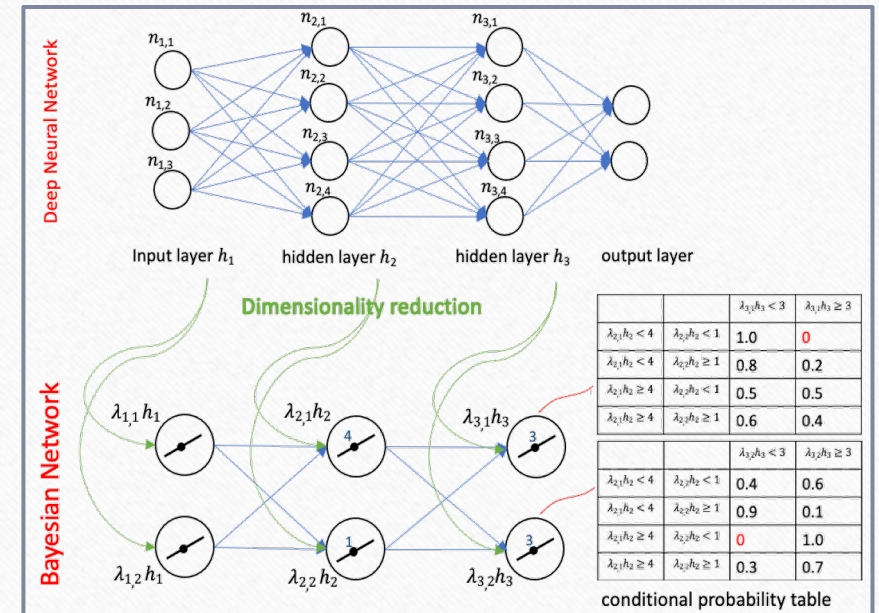1. Hidden features extraction

   Map from a high-dimensional space into a *feature space*.

2. Feature space discretisation

   Discretise each feature component into finite feature intervals.

3. Probability tables construction

   Associate each feature with a marginal or a conditional probability table.

# BN-based Latent Feature Analysis

Leverage the Bayesian Network to estimate the sensitivity of an individual feature to a controlled distribution shift.

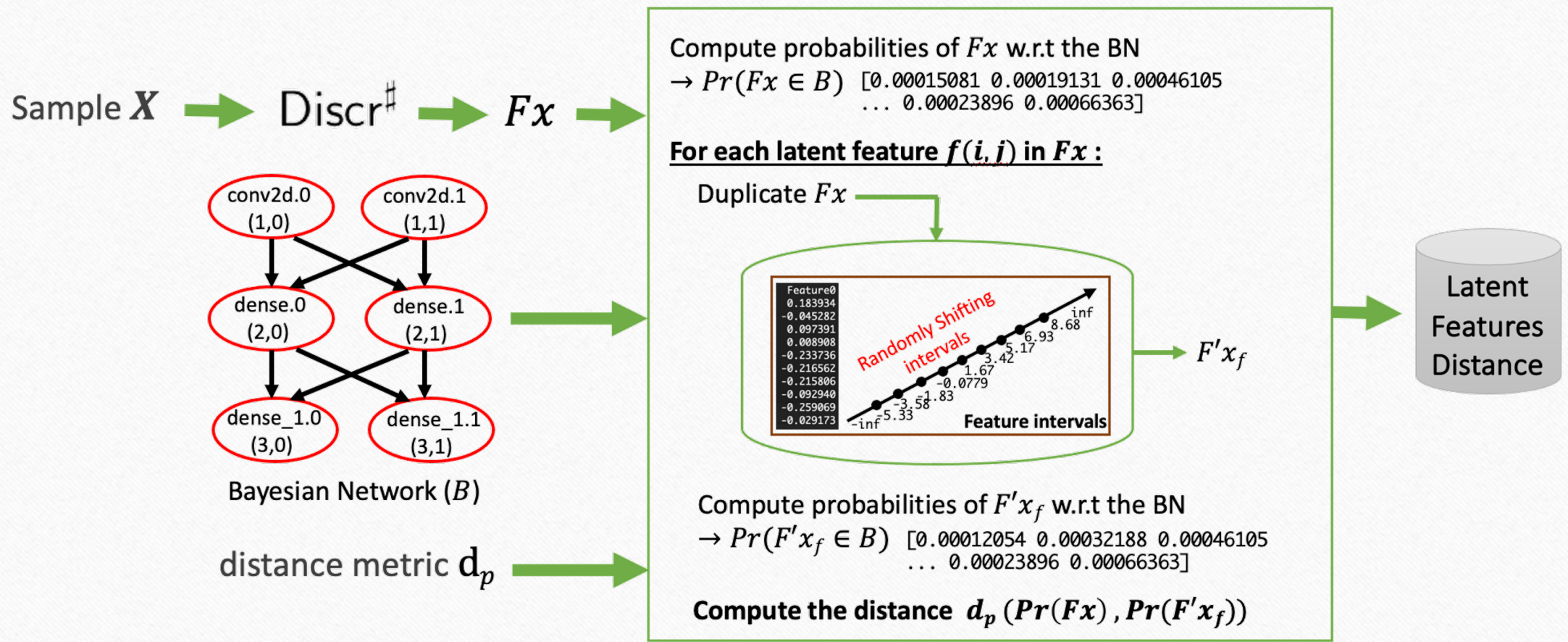1. ## Probability calculation for input sample under BN

   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.

BN analysis technique to compute the sensitivity of extracted latent features

# Experiments

## 1) Feature Importance

Associate each extracted feature $f$ with a weight $w_f$ based on the measured sensitivity distance.

**Higher distribution change → Higher importance score**

| distance<br>perturbed feature | $d_{L_1}$ | $d_{L_2}$ | $d_{L_\infty}$ | $d_{JS}$ | $d_{corr}$ | $d_{cos}$ | $d_{MSE}$ | $d_{RMSE}$ | $d_{MAE}$ | $d_{AF}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| $(1,0)$ | 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 |

Example distance measures of the MNIST model

$$w_f = \frac{e^{d_f}}{\sum_{f \in T} e^{d_f}}$$

$$w_{(1,1)} = 0.192$$
$$w_{(2,1)} = 0.182$$

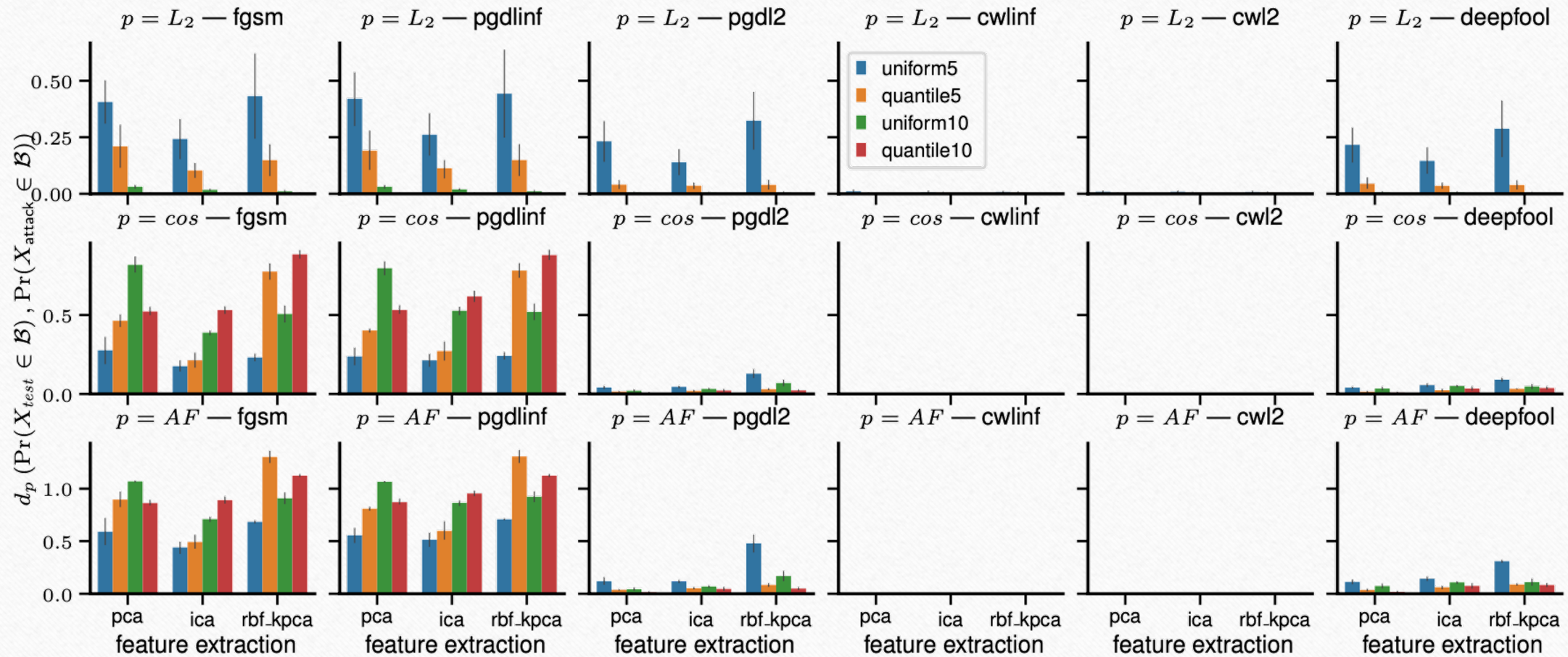UNIVERSITY OF
LIVERPOOL

# Experiments

## 2) Sensitivity to Adversarial Distribution Shift

Assess whether the BN abstraction can detect the shift in the distribution of inputs that occurs when the NN is subject to adversarial examples.

- 👽 **fgsm** is the Fast Gradient Sign Method;
- 👽 **pgdlinf** and **pgdl2** are the Projected Gradient Descent approach;
- 👽 **cwlinf** and **cwl2** ;
- 👽 **Deepfool.**

→ Generate an adversarial dataset $X_{attack}$ from the validation dataset $X_{test}$
→ Calculate the distance between their probability vectors

# Sensitivity to Adversarial Distribution Shift (MNIST)

# Results

Computing distances between two BN probability distributions clean and perturbed by intervals-shift or adversarial attacks
→ **Detect distribution shift**
→ **Reveal important features**

# Conclusion

**Advanced a novel technique that employs a BN abstraction to investigate how to measure the importance of high level features when they are used by the neural network to make classification decisions**

UNIVERSITY OF
LIVERPOOL

# Utility of the Feature Weights

First, **visualising** the most **important features** provides insight into the model's internal decisions by highlighting dominating regions in the feature space.

Second, using the **importance measurement** to design high-level **testing** metrics that evaluate the robustness of the DNN.

Third, utilising the obtained **importance** in the **training** process and force the DNN to adjust its parameters according to the most relevant features to the prediction.

UNIVERSITY OF LIVERPOOL

# THANK YOU!