# Diverse and Amortised Counterfactual Explanations for Uncertainty Estimates

Dan Ley, Umang Bhatt, and Adrian Weller

University of Cambridge

## Abstract

To interpret uncertainty estimates from differentiable probabilistic models, [1] proposed generating a single Counterfactual Latent Uncertainty Explanation (CLUE) for a given data point where the model is uncertain. [2] formulated $\delta$-CLUE, the set of CLUEs within a $\delta$ ball of the original input in latent space – however, we find that many CLUEs generated by this method are very similar, hence redundant. We propose DIVerse CLUEs ($\nabla$-CLUEs), a set of CLUEs which each provide a distinct explanation. We further introduce GLobal AMortised CLUEs (GLAM-CLUEs), which represent amortised mappings that apply to specific groups of uncertain inputs, taking them and efficiently transforming them in a single function call into inputs that a model will be certain about. Our experiments show that $\nabla$-CLUEs and GLAM-CLUEs both address shortcomings of CLUE and provide beneficial explanations of uncertainty estimates to practitioners.

## Introduction

[1] proposes a method for finding an explanation of a model's predictive uncertainty of a given input by searching in the latent space of an auxiliary deep generative model (DGM), identifying a single possible change to the input such that the model becomes more certain in its prediction. This is termed CLUE (Counterfactual Latent Uncertainty Explanation). However, there are limitations to CLUE, including the lack of a framework to deal with a potential diverse set of plausible explanations, despite proposing methods to generate them. CLUE introduces a latent variable DGM with decoder $\mu_\theta(\mathbf{x}|\mathbf{z})$ and encoder $\mu_\phi(\mathbf{z}|\mathbf{x})$. $\mathcal{H}$ refers to any differentiable uncertainty estimate of a prediction $\mathbf{y}$. CLUE minimises: $\mathcal{L}(\mathbf{z}) = \mathcal{H}(\mathbf{y}|\mu_\theta(\mathbf{x}|\mathbf{z})) + d(\mu_\theta(\mathbf{x}|\mathbf{z}), \mathbf{x}_0)$ to yield $\mathbf{x}_{\text{CLUE}} = \mu_\theta(\mathbf{x}|\mathbf{z}_{\text{CLUE}})$ where $\mathbf{z}_{\text{CLUE}} = \operatorname{argmin}_{\mathbf{z}} \mathcal{L}(\mathbf{z})$. We propose $\nabla$-CLUE and GLAM-CLUE, full details of which can be found in the paper.
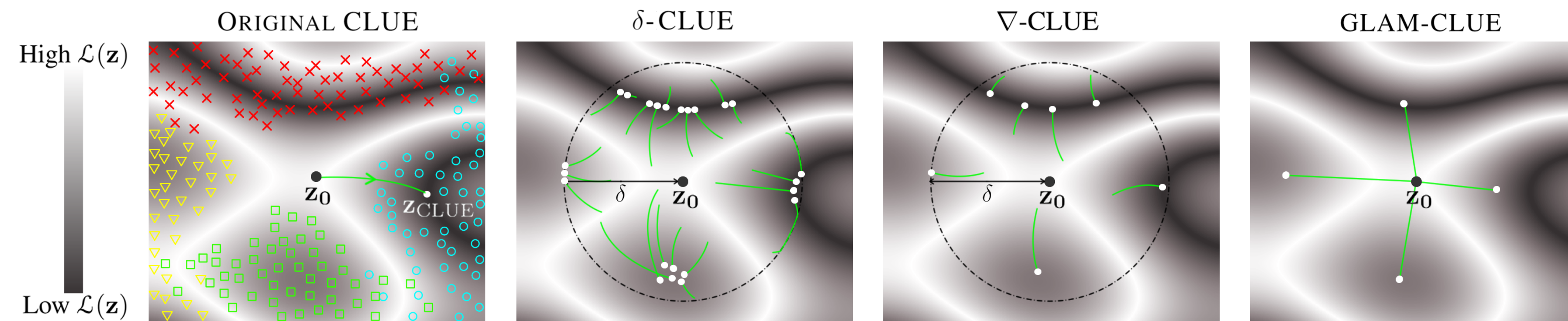


Figure 1: Conceptual colour map of objective function $\mathcal{L}(z)$ with $\mathbf{z}_0$ located in high cost region. White circles indicate explanations found. Left: Gradient descent to region of low cost [1]. Training points in colour. Left Centre: Gradient descent constrained to $\delta$-ball [2]. Diverse starting points yield diverse local minima, albeit with many redundant solutions. Right Centre: Direct optimisation for diversity ($\nabla$-CLUE). Right: Efficient mappings without gradient descent (GLAM-CLUE).

## DIVerse CLUE

$\delta$-CLUE [2] introduces a method for generating a set of CLUEs by restricting the search in latent space to a ball of radius $\delta$. However, many CLUEs found therein are redundant. We introduce metrics $D$ (Table 1) to measure the diversity in sets of CLUEs such that we can optimise for it directly: we term this DIVerse CLUE ($\nabla$-CLUE). By optimising simultaneously over $k$ counterfactuals, we minimise (note that we apply the function $D$ in latent space here):

$$\mathcal{L}(\mathbf{z}_1, ..., \mathbf{z}_k) = -\lambda_D D(\mathbf{z}_1, ..., \mathbf{z}_k) + \frac{1}{k}\sum_{i=1}^{k}\mathcal{L}(\mathbf{z}_i)$$

where $\quad \mathcal{L}(\mathbf{z}_i) = \mathcal{H}(\mathbf{y}|\mu_\theta(\mathbf{x}|\mathbf{z}_i)) + d(\mu_\theta(\mathbf{x}|\mathbf{z}_i), \mathbf{x}_0)$,

to yield $\quad X_{\text{CLUE}} = \mu_\theta(X|Z_{\text{CLUE}})$

where $\quad Z_{\text{CLUE}} = \operatorname*{argmin}_{\mathbf{z}_1, ..., \mathbf{z}_k} = \mathcal{L}(\mathbf{z}_1, ..., \mathbf{z}_k)$.

| Diversity Metric | Function ($D$) |
|---|---|
| Determinantal Point Processes | $\det(\mathbf{K})$ where $\mathbf{K}_{i,j} = \dfrac{1}{1 + d(\mathbf{x}_i, \mathbf{x}_j)}$ |
| Average Pairwise Distance | $\dfrac{1}{\binom{k}{2}}\sum_{i=1}^{k-1}\sum_{j=i+1}^{k} d(\mathbf{x}_i, \mathbf{x}_j)$ |
| Coverage | $\dfrac{1}{d'}\sum_{i=1}^{d'}\left(\max_j(\mathbf{x}_j - \mathbf{x}_0)_i + \max_j(\mathbf{x}_0 - \mathbf{x}_j)_i\right)$ |
| Prediction Coverage | $\dfrac{1}{c'}\sum_{i=1}^{c'}\max_j[(\mathbf{y}_j)_i]$ |
| Distinct Labels | $\dfrac{1}{c'}\sum_{j=1}^{c'}\mathbf{1}_{[\exists i\,:\,y_i=j]}$ |
| Entropy of Labels | $-\dfrac{1}{\log c'}\sum_{j=1}^{c'}p_j(k)\log p_j(k)$ |

Table 1: Diversity metrics $D$ with arbitrary distance metric $d$.



Figure 2: Effect of $\lambda_D$ on diversity. DPP, APD and Coverage metrics applied to the set of $k = 10$ $\nabla$-CLUEs.



| Original $\mathcal{H} = 1.13$ | Input DBM $\mathcal{H} = 0.77$ | Latent DBM $\mathcal{H} = 0.1$ | Input NN $\mathcal{H} = 0.7$ | Latent NN $\mathcal{H} = 0.42$ | GLAM 1 $\mathcal{H} = 0.6$ | GLAM 2 $\mathcal{H} = 0.39$ | CLUE $\mathcal{H} = 0.47$ |
|---|---|---|---|---|---|---|---|
| | $d = 62.9$ $\rho = 3.3$ | $d = 58.5$ $\rho = 2.0$ | $d = 42.6$ $\rho = 4.5$ | $d = 58.1$ $\rho = 3.8$ | $d = 30.5$ $\rho = 0.8$ | $d = 32.9$ $\rho = 1.0$ | $d = 26.2$ $\rho = 1.3$ |

Figure 3: Comparison of explanations for an uncertain input (left) by the baselines, GLAM-CLUE, and CLUE. $\mathcal{H}$ is uncertainty, $d$ is input space distance, $\rho$ is latent space distance. Low $\mathcal{H}$ in baselines have unrealistically high $d$ from the original.

## GLobal AMortised CLUE

We desire a computationally efficient method that only requires a finite portion of the dataset from which we learn global properties of uncertainty; we propose GLobal AMortised CLUE (GLAM-CLUE), a method that achieves this with considerable speedu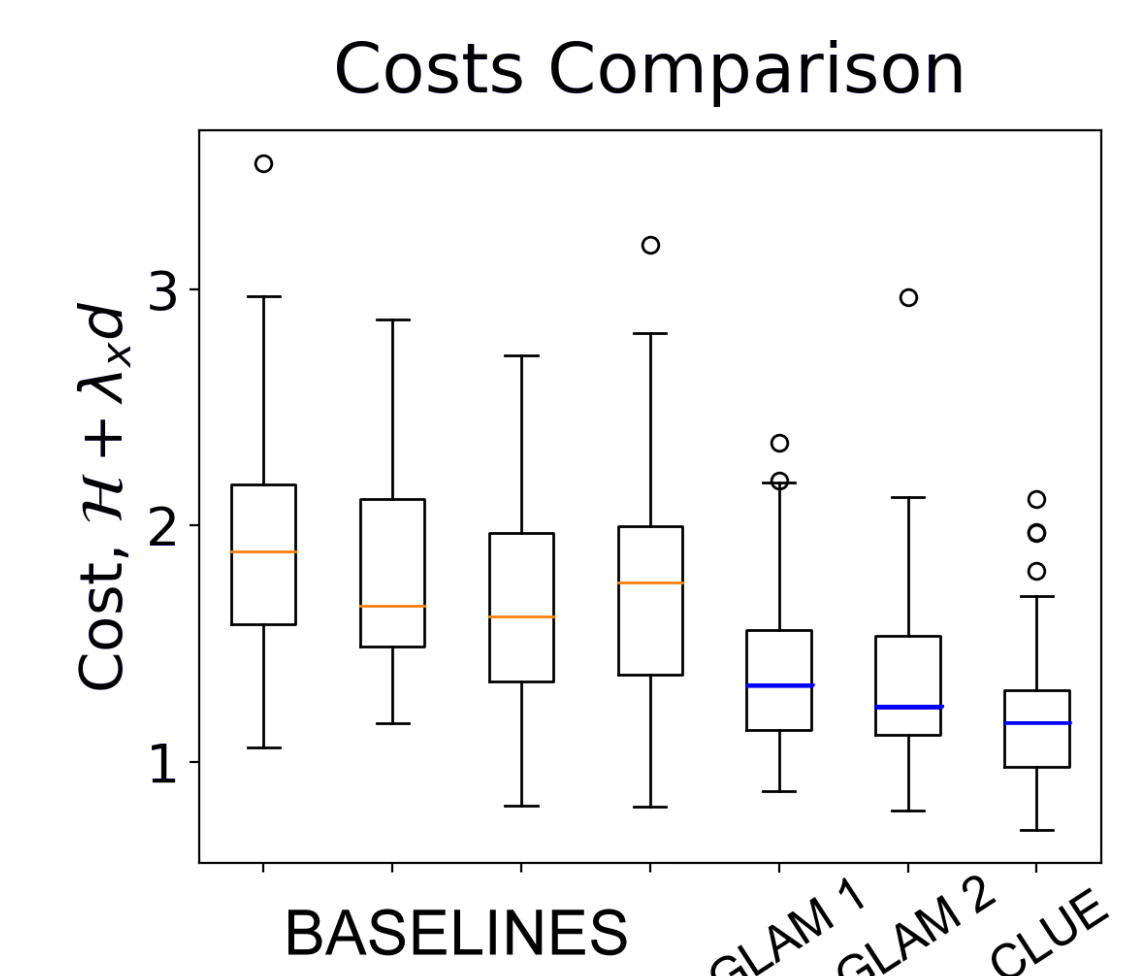ps. High certainty points are taken from the training data to learn such mappings (we call this GLAM 1), but we demonstrate improvements by instead using CLUEs generated from uncertain points in the training data (GLAM 2). At inference time, GLAM-CLUE performs significantly faster than CLUE by **average CPU time** (Table 2). We show that $\nabla$-CLUE and GLAM-CLUE address the shortcomings of CLUE.



Figure 4: GLAM-CLUE vs baselines when mapping uncertain 7s to certain 7s in MNIST. Total costs $\mathcal{H} + \lambda_x d$.

| Input DBM | Latent DBM | Input NN |
|---|---|---|
| 0.0306 | 0.0262 | 0.0236 |

| Latent NN | GLAM-CLUE | CLUE |
|---|---|---|
| 0.0245 | 0.0238 | 4.68 |

Table 2: Average CPU time in **seconds** to compute **one** MNIST counterfactual explanation.

## References

[1] Javier Antorán, Umang Bhatt, Tameem Adel, Adrian Weller, and José Miguel Hernández-Lobato. Getting a CLUE: A method for explaining uncertainty estimates. In *International Conference on Learning Representations*, 2021.

[2] Dan Ley, Umang Bhatt, and Adrian Weller. $\delta$-CLUE: Diverse sets of explanations for uncertainty estimates. In *ICLR Workshop on Security and Safety in Machine Learning Systems*, 2021.

## Acknowledgements