

000

001

002

003

004

005

006

007

008

009

010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

054

# Using the Path of Least Resistance to Explain Deep Networks

Anonymous Authors<sup>1</sup>

## Abstract

Integrated Gradients (IG), a widely used path-based attribution method, assigns importance scores to input features by integrating gradients of the models along a straight path from a baseline to the input. While effective in certain cases, we show that choosing straight paths can lead to flawed attributions. In this paper, we identify how these misattributions arise. As a solution, we propose a new approach that considers the input space as a Riemannian manifold and computes attributions by integrating gradients of the model along geodesics. We call our approach *Geodesic Integrated Gradients*. To approximate geodesic paths, we introduce two methods: a  $k$ -nearest neighbour-based approach for simpler models and a Stochastic Variational Inference-based method for more complex ones. Furthermore, in our experiments with both synthetic and real world data, we demonstrate that our approach outperforms existing explainability methods including the original IG.

## 1. Introduction

The use of deep learning models has risen in many applications. With it, so too has the desire to understand why these models make certain predictions. These models are often referred to as “opaque”, as it is difficult to discern the reasoning behind their predictions (Marcus, 2018). Additionally, deep learning models can inadvertently learn and perpetuate biases found in their training data (Sap et al., 2019). To create fair and trustworthy algorithms, it is essential to be able to explain a model’s output (Das & Rad, 2020).

Some examples of the methods proposed to explain neural networks include Gradient SHAP (Lundberg & Lee, 2017), Integrated Gradients (Sundararajan et al., 2017) and (Kapishnikov et al., 2021).

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

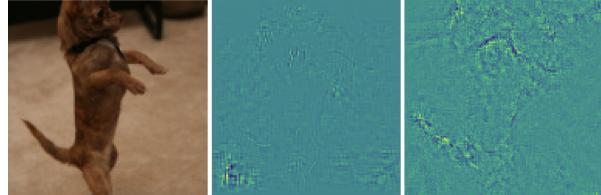


Figure 1. Comparing attributions of Integrated Gradients (middle image) with Geodesic Integrated Gradient (right image). In this work we show how to generate sharper explanations by integrating gradients along the geodesic line, as opposed to the straight line.

Significant effort has been dedicated to designing explanation methods that satisfy certain desirable axioms. This is due to the lack of ground truth for evaluating them. The axioms can ensure that the explanations are principled. One of the most successful axiomatic methods is Integrated Gradients (IG) (Sundararajan et al., 2017). Consider a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , representing the neural network and an input vector  $\mathbf{x} \in \mathbb{R}^n$ . Furthermore, consider a baseline input vector  $\bar{\mathbf{x}} \in \mathbb{R}^n$  (typically chosen such that the network gives baseline a near zero score). IG explains the network by quantifying how much of the difference  $f(\mathbf{x}) - f(\bar{\mathbf{x}})$  can be attributed to the  $i$ th dimension of  $\mathbf{x}$ ,  $\mathbf{x}_i$ .

Integrated Gradient gives attribution  $IG_i$  to the  $i$ th dimension of the input by solving the following path integral

$$IG_i(\mathbf{x}) = (\mathbf{x}_i - \bar{\mathbf{x}}_i) \int_0^1 \frac{\partial f(\gamma(t))}{\partial \mathbf{x}_i} dt, \quad (1)$$

where  $\gamma(t) = \bar{\mathbf{x}} + t(\mathbf{x} - \bar{\mathbf{x}})$  is a straight path from the baseline to input. The claim of the creators of IG is that Eq. 1 tells us how the model got from predicting essentially nothing at  $\bar{\mathbf{x}}$  to giving the prediction at  $\mathbf{x}$ . Considering gradients represent the rate of change of functions, the above expression should tell us how scaling each feature along the path affects the increase in the network score for the predicted class.

In this paper, we demonstrate that defining attributions along straight paths in Euclidean space can lead to misattributions. We examine the consequences of these issues through examples in computer vision, as shown in Fig. 1, alongside simpler illustrative cases in Fig. 2. To address these challenges, we introduce **Geodesic Integrated Gradients**<sup>1</sup>, a gener-

<sup>1</sup>The code for the attribution methods and reproducing the

alisation of IG that replaces straight paths with geodesic ones. These geodesics are defined on a Riemannian manifold, characterised by the model’s input space and a metric induced by the model’s gradients. This approach mitigates the identified pitfalls while retaining all the axioms of IG.

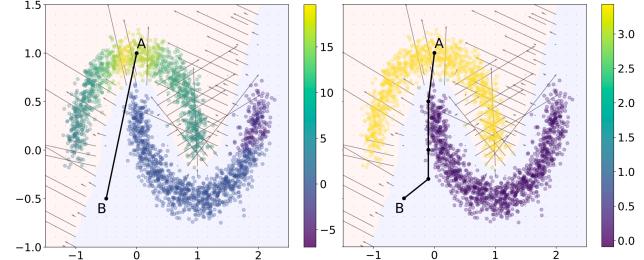
Before making the case for our Geodesic Integrated Gradient, let us first show an example of an artefact that can arise from choosing straight paths, generating explanations which do not reflect the true behaviour of a model.

We highlight this issue on a half-moons classification task. We train a simple multi-layer perceptron (MLP) with 3 layers, ReLU activations and a cross-entropy loss to distinguish the upper moon from the lower one. The cross-entropy is split into a final log-softmax activation and a negative log-likelihood loss, so that we can explain probabilities.

We now compute Integrated Gradients, Eq. 1, for this model on the test data. Let us consider a baseline and input pair, such that the baseline is outside of either half-moon, for example at (-0.5, -0.5). This is a good choice of baseline, since network should assign near-zero score to it. Let us call the feature in the vertical axis the 1st component of  $x$ . In Fig. 2 we illustrate the attribution of this feature,  $IG_1(\mathbf{x})$ , for each point using the colour map. One should expect to see all the points sufficiently above the decision boundary to receive equally high attributions. Intuitively, this is expected, because if a point is above the decision boundary, its  $x_1$  component is an important factor in the classification. However, for a model that is very skillful at the classification task, since the point is significantly above the decision boundary, going slightly down should not make any difference. Because in such a model’s score should not change significantly anywhere other than near the decision boundary. However, we can see in Fig. 2 that some points on the upper moon receive much higher  $x_1$ -attribution than others. These are points such that a straight line from the baseline to them mostly falls on high gradient regions. This does not reflect the model’s behaviour. A similar point could be made about the horizontal axis. This is in contrast with Fig. 2, where we show our method gives equally high attribution to all points sufficiently above the decision boundary, with different shades for the points closer to the boundary in  $x_1$  direction. In Section 3.1 we present the details of how our method that achieves the results presented in this figure.

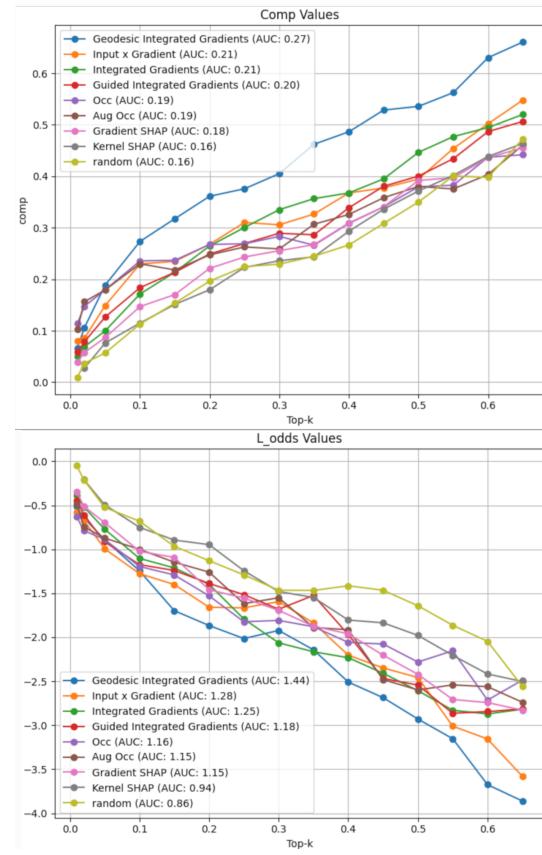
Before giving the formal method in section 2, let us discuss the intuition behind our Geodesic IG. We want the path in Eq. 1 to be such that it avoids regions with high model gradient. This is because failing to do so would superficially increase the result of the integral, leading to the types of

experiments is provided on <https://github.com/sina-salek/geodesic-ig>



**Figure 2. Integrated Gradients attributions.** The colour map represents  $IG_1(\mathbf{x})$  for each point  $\mathbf{x}$ , with  $(-0.5, -0.5)$  (point B) as the baseline. Points around A have much higher attributions than other points on the top moon, despite all being sufficiently above the decision boundary of the MLP. This is due to the path being close to the decision boundary, resulting in high gradients (gray arrows) along this path.

artefacts illustrated in Fig. 2. Therefore, we should try to find the path of least resistance, as this is the path that avoids steep gradients as much as possible. As we shall see in section 2, the input space can be viewed as a Riemannian manifold with a metric derived from the model gradients.



**Figure 3. Metrics comparison** Evaluation of our results via comprehensiveness (De Young et al., 2019) and log odds (Shrikumar et al., 2017) show Geodesic Integrated Gradients significantly outperforming other methods.

110 The path of least resistance between a chosen baseline and  
 111 input, therefore, is the geodesic path between the two points.

112 In section 2, we develop two methods to approximate the  
 113 geodesic path between two points on the manifold. We use  
 114 the first method,  $k$ NN-based for simpler manifolds, and the  
 115 second one, based on Stochastic Variational Inference, for  
 116 more complex manifolds. We then show that geodesic IG  
 117 satisfies all the axioms of IG.

118 In Section 3, we demonstrate the effectiveness of the  
 119 Geodesic IG method on the real-world Pascal VOC 2012  
 120 dataset (Everingham et al.). Our results outperform existing  
 121 methods, as we evaluate using various metrics. We preview  
 122 the results of this experiment in Fig. 3

123 Section 4 reviews related work, including the comparison of  
 124 Geodesic IG with other methods that attempt to overcome  
 125 the shortcomings of Integrated Gradients.

## 128 2. Method

130 In section 1, we gave the intuition that using geodesic paths  
 131 can correct the misattribution in IG that arise from integrating  
 132 along straight paths. Let us now formalise this idea.

### 134 2.1. Geodesic distance formulation.

136 Let us define a neural network as a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  
 137 where  $n$  is the dimension of the input space. Let us also  
 138 define  $\mathbf{x}$  a point in this input space. We denote the Jacobian  
 139 of  $f$  at  $\mathbf{x}$  as  $J_{\mathbf{x}}$ .

141 Using Taylor’s theorem, for a vector  $\delta$  with an infinitesimal  
 142 norm:  $\forall \epsilon, \|\delta\| \leq \epsilon$ , we have:

$$145 \|f(\mathbf{x} + \delta) - f(\mathbf{x})\| \approx \|J_{\mathbf{x}}\delta\| \approx \delta^T J_{\mathbf{x}}^T J_{\mathbf{x}} \delta \quad (2)$$

147 Using equation 2, we can now define a tangent space  $T_{\mathbf{x}}M$   
 148 of all  $\delta$ , equipped with a local inner product  $G_{\mathbf{x}}$ :

$$151 \langle \delta, \delta' \rangle_{\mathbf{x}} = \delta^T G_{\mathbf{x}} \delta' = \delta^T J_{\mathbf{x}}^T J_{\mathbf{x}} \delta' \quad (3)$$

153 As a result, we can view the input space as a Riemannian  
 154 manifold  $(\mathbb{R}^n, G)$ , where the Riemannian metric  $G$   
 155 is defined above. On this manifold, the length of a curve  
 156  $\gamma(t) : [0, 1] \rightarrow \mathbb{R}^n$  is defined as:

$$159 L(\gamma) = \int_0^1 \sqrt{\langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{\gamma(t)}} dt \\ 160 = \int_0^1 \|\partial_t f(\gamma(t)) \times \dot{\gamma}(t)\| dt, \quad (4)$$

where  $\dot{\gamma}(t)$  is the derivative of  $\gamma(t)$  with respect to  $t$ . The  
**geodesic distance**, denoted  $L^*$ , between  $\mathbf{a}$  and  $\mathbf{b}$  is then  
 defined as the minimum length among curves  $\gamma$  such that  
 $\gamma(0) = \mathbf{a}$  and  $\gamma(1) = \mathbf{b}$ . We also call **geodesic path**  
 the curve  $\gamma^*$  which minimises the length  $L$ . This path can  
 be interpreted as the shortest path between  $\mathbf{a}$  and  $\mathbf{b}$  in the  
 manifold.

*Remark 2.1.* We can infer from Equation 4 that the geodesic  
 path avoids as much as possible high-gradients regions. This  
 is the main desired property of a path to be used for path-  
 based attributions. Representing the path of least resistance,  
 the geodesic path circumvents superficially high values of  
 attributions.

### 2.2. Approximation of the geodesic with $K$ Nearest Neighbours.

Computing the exact geodesic would require computing  $L$   
 on an infinite number of paths  $\gamma$ , which is not possible in  
 practice. However, several methods have been proposed  
 to approximate this value. We draw from previous work  
 (Yang et al., 2018; Chen et al., 2019) and present one with  
 desirable characteristics.

First, we compute the K Nearest Neighbours (kNN) algorithm  
 on points between (and including) input and baseline. These points can be either sampled or generated. The  
 geodesic distance between two neighbouring points,  $\mathbf{x}_i$  and  
 $\mathbf{x}_j$ , can be approximated by a straight path  $\mathbf{x}_i + t \times (\mathbf{x}_j - \mathbf{x}_i)$ . We have the above approximation because for dense enough  
 data, the euclidean distance between neighbouring points  
 is a good approximation of the geodesic distance. This reflects the fact that a small region of a Riemannian manifold,  
 called Riemann neighbourhood, is locally isometric to a Euclidean space<sup>2</sup>. So the geodesic distance between the two  
 neighbouring points is approximated by:

$$L_{ij}^* = \int_0^1 \|\partial_t f(\mathbf{x}_i + t \times (\mathbf{x}_j - \mathbf{x}_i)) \times (\mathbf{x}_i - \mathbf{x}_j)\| dt \\ = \|\mathbf{x}_i - \mathbf{x}_j\| \int_0^1 \|\partial_t f(\mathbf{x}_i + t \times (\mathbf{x}_j - \mathbf{x}_i))\| dt \quad (5)$$

Equation 5 corresponds to the original Integrated Gradients  
 method, albeit with the norm. This integral can be approximated  
 by a Riemannian sum similarly to (Sundararajan et al., 2017):

$$L_{ij}^* \approx \|\mathbf{x}_i - \mathbf{x}_j\| \sum_{k=0}^m \|\partial f(\mathbf{x}_i + \frac{k}{m} \times (\mathbf{x}_j - \mathbf{x}_i))\| \quad (6)$$

For input-baseline pair,  $\mathbf{x}$  and  $\bar{\mathbf{x}}$ , we can now see the set

<sup>2</sup>We shall further formalise this intuition later in this section.

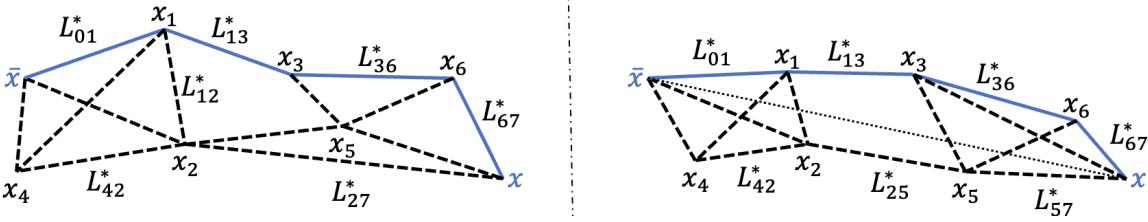


Figure 4. **Method overview.** For an input  $\mathbf{x}$ , a baseline  $\bar{\mathbf{x}}$ , and a set of points  $\mathbf{x}_i$ , we compute the kNN graph using the euclidean distance (dashed lines). For each couple  $(\mathbf{x}_i, \mathbf{x}_j)$ , we then compute the integrated gradients  $L_{ij}^*$  using Equation 6. For clarity, not all  $L_{ij}^*$  are present on the figure. 0 and 7 represent  $\bar{\mathbf{x}}$  and  $\mathbf{x}$  respectively. Using the resulting undirected weighted graph, we use the Dijkstra algorithm to find the shortest path between  $\mathbf{x}$  and  $\bar{\mathbf{x}}$  (blue continuous lines). On the left, the points  $\mathbf{x}_i$  are provided while, on the right, the points are generated along the straight line between  $\mathbf{x}$  and  $\bar{\mathbf{x}}$  (dotted line).

$(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{x}_i)$  as a weighted graph, with the weights being the geodesic distances between two neighbors  $L_{ij}^*$ . To compute the geodesic path between  $\mathbf{x}$  and  $\bar{\mathbf{x}}$ , we can use a shortest path algorithm, such as Dijkstra or A\* with the euclidean distance as the heuristic.

The resulting Geodesic Integrated Gradients corresponds to the sum of the gradients along this shortest path:

$$\text{Geodesic IG}_i(\mathbf{x}) = (x_i - \bar{x}_i) \sum_{k=0}^m \int_0^1 \frac{\partial f(\mathbf{x}^k + t \times (\mathbf{x}^{k+1} - \mathbf{x}_k))}{x_i^k} dt \quad (7)$$

where  $\mathbf{x}^k$  are the points along the shortest path. The integrals in Equation 7 can also be approximated with Riemannian sums.

The gradients between each pair of neighbours can also be estimated in batches to speed up the attribution computation. Moreover, several inputs' attributions can be computed together, with similar speed as IG: if we want to compute the attribution of  $N$  inputs, with 10 interpolation steps and 5 nearest neighbors, the number of gradients to calculate is  $10 \times 5 \times N = 50N$ , which amounts to computing IG with 50 steps. This does not include the computation of the shortest path, which is for instance  $O(N^2)$  for Dijkstra algorithm. Please also refer to Figure 4 for an illustration of this method.

**Assumption of the approximation.** Here we formalise the intuition that, for a pair of neighbours, the geodesic path between them is close to the euclidean one. Notice that the derivative of the neural network  $f$  is Lipschitz continuous,

$$\exists K \forall \mathbf{x}, \mathbf{y}, \|J_x - J_y\| \leq K \times \|\mathbf{x} - \mathbf{y}\|. \quad (8)$$

Equation 8 is equivalent to the Hessian of  $f$  being bounded. Under this assumption, if two points  $\mathbf{x}$  and  $\mathbf{y}$  are close

enough, the Jacobian of one point is approximately equal to the other: if  $\|\mathbf{x} - \mathbf{y}\| \leq \epsilon$ , then  $J_x \approx J_y$ . As a result, the length between  $\mathbf{x}$  and  $\mathbf{y}$ , for a curve  $\gamma$ , is:  $L(\gamma) \approx \int_\gamma \|J_x\| dx \approx \|J_x\| \int_\gamma dx$ . Due to the triangular inequality, the shortest path  $\gamma^*$  is then a straight line, and we have:  $L^*(\mathbf{x}, \mathbf{y}) \approx \|J_x\| \times \|\mathbf{x} - \mathbf{y}\|$ .

As a result, under this assumption, if two points are close, the geodesic path can be approximated with a straight line. Note that even though we take the path between two neighbouring points to be a straight line, we do not assume that the Jacobian of the function between the two points is constant.

**Handling disconnected graphs** An issue with the graph computed with the kNN algorithm is that it could be disconnected, in which case it could be impossible to compute a path between an input and a baseline. To alleviate this issue, we add so called “bridges” to the graph, as following: for each disconnected component, we add one link between them, specifically between two points of each component having the lowest euclidean distance. An illustration of this method is displayed on Figure 5.

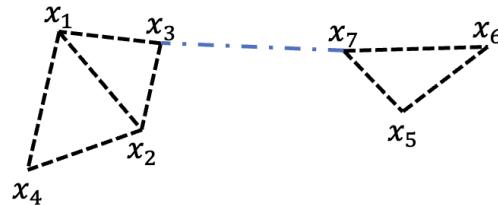


Figure 5. When the kNN graph is disconnected, as illustrated here, it would be impossible to compute Geodesic IG between certain points, for instance  $\mathbf{x}_1$  and  $\mathbf{x}_5$  here. To solve this, we add a single link between disconnected graphs, here between  $\mathbf{x}_3$  and  $\mathbf{x}_7$ .

However, we stress that this solution is not optimal, and argue that a better way of handling this issue would be to

220 avoid disconnected graphs in the first place. This can be  
 221 done by increasing the number of neighbours k.  
 222

### 223 2.3. Approximation of the geodesic with energy-based 224 sampling.

225 While our  $k$ NN-based method is effective for explaining  
 226 simpler models, its applicability diminishes as model  
 227 complexity increases. In such cases, a prohibitively large number  
 228 of samples is required between the baseline and the input to  
 229 provide accurate estimates of the geodesic path. Even with  
 230 relatively large number of samples, it is not trivial where on  
 231 the manifold to sample the points to adequately capture the  
 232 gradient landscape. Furthermore, once the points are sam-  
 233 pled, searching the graph for the shortest path will be com-  
 234 putationally too intensive. For such use-cases, we devise  
 235 another approximation procedure based on an energy-base  
 236 sampling method, such as Stochastic Variational Inference  
 237 (SVI).  
 238

239 As we observed in Remark 2.1, we would like to sample  
 240 from the shortest paths between two points on the input  
 241 space that avoids high gradient regions as much as possible.  
 242 Therefore, we want to deviate from the straight line to avoid  
 243 high gradient regions. This process can be approximated  
 244 in the following way. First we start with a straight line  
 245 between the two points and define a potential energy as  
 246 a combination of two terms. One for the distance of the  
 247 path to be optimised to the straight line, and the other is  
 248 a curvature penalty term. Finding the minimum energy  
 249 path, therefore samples approximately on the geodesic lines.  
 250 More formally, let us define the distance term:  $d(x, y) :=$   
 251  $|x - y|_2$ , and the curvature term:  $c(x) := \|\nabla f(x)\|_2$  where  
 252  $f$  is the neural network. The total energy being minimised  
 253 is

$$254 E(\gamma) = \sum_{i=1}^n \|\gamma_i - \gamma_i^0\|_2 - \beta \|\nabla f(\gamma_i)\|_2, \quad (9)$$

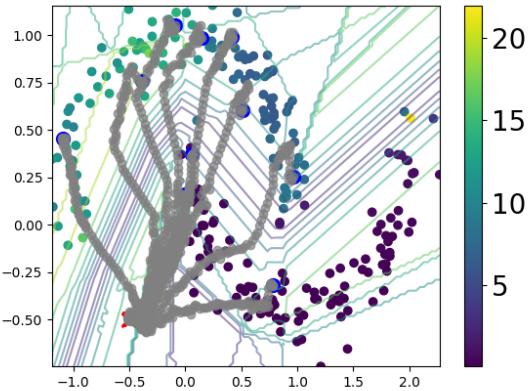
255 where  $\gamma$  is the path  $\gamma^0$  is the initial path beta controls trade-  
 256 off between distance and curvature.  
 257

258 With this energy function, one can use a suitable sampling  
 259 method, such as SVI or Hamiltonian Monte Carlo to sam-  
 260 ple points on the geodesic paths. Here we briefly describe  
 261 the SVI optimisation, as this has a good balance of compu-  
 262 tational efficiency and accuracy.  
 263

264 SVI provides a probabilistic framework for optimising paths  
 265 between input and baseline points. To achieve this, SVI  
 266 defines a probability distribution  $p(\gamma|\gamma_0)$  proportional to  
 267  $\exp(-E(\gamma))$ , where  $E(\gamma)$  is our defined potential energy.  
 268 Rather than directly sampling from this complex distribution,  
 269 we introduce a simpler variational distribution  $q(\gamma)$   
 270 parameterised by learnable means and scales. This guide  
 271 distribution takes the form of a factorised normal distribu-  
 272 tion  $\prod_i N(\mu_i, \sigma_i)$  over path deviations.  
 273

274 The optimisation proceeds by minimising the KL divergence  
 275 between  $q(\gamma)$  and the true posterior through maximisation of  
 276 the Evidence Lower Bound (ELBO). Critically, this allows  
 277 us to learn optimal parameters for  $q(\gamma)$  through gradient-  
 278 based optimisation. The learned means  $\mu_i$  define the optimal  
 279 path deviations from the initial straight-line path, while  
 280 the scales  $\sigma$  capture uncertainty in these deviations. This  
 281 probabilistic approach naturally samples of the low-energy  
 282 regions.

283 We use this method in our computer vision experiments  
 284 and show its efficacy in section . However here we find it  
 285 instructive to visualise these paths on the 2D example of  
 286 half-moons, Fig. 6, even though for simpler use-cases such  
 287 as this we would prefer the  $k$ NN method, as they are easier  
 288 to control.



289 **Figure 6. Visualising 10 random paths.** For a simple case of half-  
 290 moons with 1000 samples, we show the sampled paths between  
 291 10 pairs of points. As we see on low gradient regions the sampler  
 292 prefers straight lines, whereas in high gradient regions the path  
 293 becomes closer to perpendicular to the large gradient vectors to  
 294 cross the region as quickly as possible.

### 295 2.4. Axiomatic properties

296 The family of generalisations of Integrated Gradients to  
 297 non-straight paths, such as Eq. 7, is called *path methods* of  
 298 explanation. We see in Sundararajan et al. (2017) that all  
 299 path methods satisfy all of the axioms that IG is based on,  
 300 apart from the symmetry axiom. Therefore, here we focus  
 301 on this axiom only.

**Symmetry preserving of Geodesic IG** The symmetry  
 302 axiom is defined in the following way.

**Definition 2.2.** Consider an input-baseline pair  $\mathbf{x}$  and  $\bar{\mathbf{x}}$ ,

and a function  $f$  that is symmetric in dimensions  $i$  and  $j$ . If  $\mathbf{x}_i = \mathbf{x}_j$  and  $\bar{\mathbf{x}}_i = \bar{\mathbf{x}}_j$ , then an attribution method is Symmetry-Preserving if  $\text{attr}_i(\mathbf{x}; f) = \text{attr}_j(\mathbf{x}; f)$ , where  $\text{attr}_n(\mathbf{x}; f)$  is the attribution of  $\mathbf{x}_n$ .

(Sundararajan et al., 2017, Theorem 1) shows that IG is the only path method that satisfies symmetry for any function. However, as noted in (Kapishnikov et al., 2021), while straight path is the only one satisfying symmetry for any function, given a specific function, it may be possible to find other paths that satisfy Symmetry. We generalise this theorem for Geodesic IG on Riemannian manifolds and how that our method satisfies Symmetry for the neural network function we sample the paths for.

Here we demonstrate that Geodesic Integrated Gradients satisfies symmetry property on Riemannian manifolds induced by the input space and the gradients of the model.

Let the  $i$ th and  $j$ th dimensions of  $\gamma(t)$  be  $\gamma_i(t)$  and  $\gamma_j(t)$  respectively and  $f$  be a function differentiable almost everywhere on  $t$ . Furthermore, take  $f$  to be symmetric with respect to  $x_i$  and  $x_j$ . If  $\gamma_i(t) = \gamma_j(t)$  for all  $t \in [0, 1]$ , then we have

$$\|\partial_t f(\gamma_i(t)) \times \dot{\gamma}_i(t)\| = \|\partial_t f(\gamma_j(t)) \times \dot{\gamma}_j(t)\|, \quad (10)$$

almost everywhere on  $t$ . Therefore, the  $i$ th and  $j$ th components of Eq. 4 are equal. Furthermore, since Eq. 7 integrates along the path that is an approximation of Eq. 4, we have Geodesic  $\text{IG}_i = \text{Geodesic IG}_j$ . Indeed our geodesic paths satisfy  $\gamma_i(t) = \gamma_j(t)$  for all  $t \in [0, 1]$  on the Riemannian manifolds. To see this, let us select a baseline  $\bar{\mathbf{x}}$  and  $U$  a Riemann neighbourhood centred at  $\bar{\mathbf{x}}$ . Let us also define the geodesic path  $\gamma$  such as  $\gamma(0) = \bar{\mathbf{x}}$ . Further, define  $\mathbf{v}(t) := \gamma'(t)$ , where  $\gamma'$  is the derivative of  $\gamma$ . Then, in the local coordinates system of the neighbourhood of any point, called normal coordinates, we have  $\gamma(t) = (tv_1(t), \dots, tv_n(t))$ . Since the function is symmetric in the  $i$ th and  $j$ th dimensions, we have  $v_i$  and  $v_j$  are the same everywhere. From this, we can see that  $\gamma_i(t) = \gamma_j(t)$  for all  $t \in [0, 1]$  and therefore Geodesic IG satisfies symmetry.

### 3. Experiments

To validate our method, we performed experiments on two datasets: one is the synthetic half-moons dataset, and the other is the real-world Pascal VOC 2012 dataset.

#### 3.1. Experiments on the half-moons dataset

Here we give the details of the half-moons experiment discussed in the Introduction section. We use the half-moons dataset provided by Scikit learn (Pedregosa et al., 2011) to generate 10,000 points with a Gaussian noise of  $\mathcal{N}(0, 0.2)$ . The dataset is split into 8,000 training points and 2,000

testing ones. The model used is an MLP.

We measure two indicators of performance for each attribution method: a lack of artifacts not reflecting the model's behaviour, and a low variation of attributions between close points. To that goal, we use two metrics, **purity** and **standard deviation**, defined in the following way. We know that a well-trained model should classify about half of the data points as upper moon, class 1, and the other half as lower moon, class 0. Such a model should consider both features of each point important for the classification into class 1. Therefore, for such a model, we expect in a good attribution method, the top 50% points as ranked by the quantity  $\widetilde{\text{attr}}(\mathbf{x}; f) = \sum_{i=0}^1 |\text{attr}_i(\mathbf{x}; f)|$ , to be classified as 1 (assuming the baseline is chosen as a point to which the network gives a near-zero score). With this in mind, purity is defined as

$$\text{Purity} = \frac{1}{N/2} \sum_{\mathbf{x}, \widetilde{\text{attr}}(\mathbf{x}; f) \in \text{Top 50\% of all attr}} \text{argmax}(f(\mathbf{x})), \quad (11)$$

where  $N$  is the number of data points. We see that this is the average value of the predicted class labels for half of the points. From the above, we infer that, for a well-trained model, we prefer an attribution method that results in the purity close to 1. In contrast, a random attribution method in this case would result in the purity score of 0.5.

For the second metric, standard deviation, points that belong to the same moon to have similar  $\widetilde{\text{attr}}(\mathbf{x}; f)$ . Therefore, define

$$\text{Std}_i = \text{std}\{\widetilde{\text{attr}}(\mathbf{x}; f), \text{argmax}(f(\mathbf{x})) \in \text{Moon } i\}. \quad (12)$$

We expect this standard deviation metric to be low for the points belonging to either class.

In this experiment, we compare the results of attributions from Geodesic IG with the original IG, as well as more recent comparable methods including Enhanced IG (Jha et al., 2020), GradientShap (Lundberg & Lee, 2017), SmoothGrad (Smilkov et al., 2017).

For all of the methods, we use  $(-0.5, -0.5)$  as a baseline. Enhanced IG is an improvement over IG where the kNN algorithm is used to avoid computing gradients on paths on out of sample distributions. The chosen number of neighbours for the kNN part of both Enhanced IG and Geodesic IG is 5. We perform an ablation study of this parameter in Appendix A.

The quantitative analysis of the results are given in Table 1, where we see that Geodesic IG significantly outperforms all other methods on all metrics. In particular, notice that all other methods perform relatively poorly on  $\text{STD}_1$ . This is because for this metric, the integration path has to cross the

Geodesic Integrated Gradients

Method	Purity $\uparrow$	$STD_0 \downarrow$	$STD_1 \downarrow$
GradientShap	0.761	1.90	5.01
IG	0.802	1.44	4.915
SmoothGrad	0.800	1.42	4.720
Enhanced IG	0.738	1.23	1.876
<b>Geodesic IG</b>	<b>0.978</b>	<b>0.237</b>	<b>0.267</b>

Table 1. Evaluation of different attribution methods on a half-moons dataset with Gaussian noise  $\mathcal{N}(0, 0.2)$ . The results over 5 different seeds are averaged, with the corresponding standard deviation in brackets. We present in Appendix A more results with different amounts of Gaussian noise.

decision boundary. However, the other methods might cross the decision boundary differently and regardless of the function’s gradient, resulting in spuriously different attributions for different points belonging to class 1. In Appendix A we see that the gap between the performance of geodesic IG and the other methods increases as the Gaussian noise of the half-moons increases. To provide better understanding of these results we present more analysis on this dataset in Section 4.

### 3.2. Experiments on the Pascal VOC 2012 dataset

Now we would like to test our method on a real-world dataset. To this aim, we use the Pascal VOC 2012 dataset, which consists of labelled images (Everingham et al.). We trained a classification head on this dataset and added it to the pre-trained ConvNext model from torch vision to generate predictions to be explained (He et al., 2016). We also perform the analysis on 100 randomly sampled images from the test set.

We compare our method with various baselines: Integrated Gradients, GradientShap, InputXGradients (Shrikumar et al., 2016), KernelShap (Lundberg & Lee, 2017), Occlusion (Zeiler & Fergus, 2014) Augmented Occlusion (Tonekaboni et al., 2020) and Guided IG (Kapishnikov et al., 2021).

Where an explainer needs a baseline, such as in IG as well as our method, a uniformly black image was chosen as baseline.

To evaluate the performance of an attribution method, we use 2 different metrics:

- **Comprehensiveness** (DeYoung et al., 2019): We mask the top  $k\%$  most important features in absolute value, and compute the average change of the predicted class probability compared with the original image. A higher score is better as it indicates masking these features

Method	AUC-Comp $\uparrow$	AUC-LO $\uparrow$
Input X Gradients	0.21	1.28
GradientShap	0.18	1.15
IG	0.21	1.25
Random	0.16	0.86
Kernel Shap	0.16	0.94
Occlusion	0.19	1.16
Aug Occlusion	0.19	1.15
Guided IG	0.20	1.18
<b>Geodesic IG</b>	<b>0.27</b>	<b>1.44</b>

Table 2. Evaluation of different attribution methods on 100 randomly sampled images from the Pascal VOC test set. Fig. 3 shows the curves where these metrics are extracted from.

results in a large change of predictions.

- **Log-odds**<sup>3</sup> (Shrikumar et al., 2017): We mask the top  $k\%$  most important features in absolute value, and measure the negative logarithmic probabilities on the predicted class compared with the original one. Lower scores are better.

We calculated these two metrics for a range of top- $k\%$ , from 1% to 56%, which is shown in Fig. 3. To summarise the values from different top- $k\%$ , we calculate the area between the curves and the horizontal line at  $y = 0$ , a.k.a. area under the curve, and report them in Table 2. Fig. 1 provides a qualitative comparison between Geodesic IG and Integrated Gradients. The analysis shows that Geodesic IG outperforms other methods in explaining the model’s behaviour on the dataset, and it does so with a striking gap in the case of comprehensiveness.

The one downside of using Geodesic IG for explaining complex deep learning models is the computational cost. Running the above experiment on 100 images took 23 hours using an L4 GPU. However, this would be a reasonable price to pay to get more accurate explanations in use-cases where understanding the model is more critical.

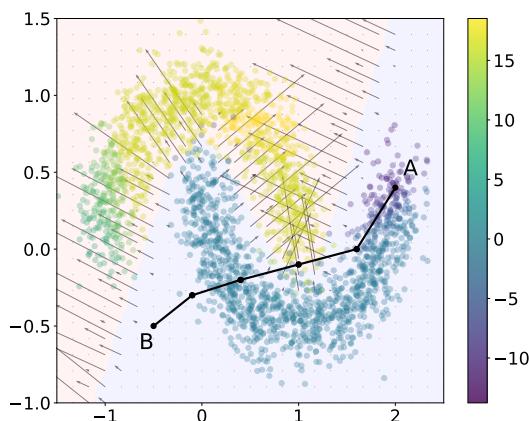
## 4. Related Work

Approximating geodesic paths is a widely studied area of research, and many methods to do so have been developed. For a comprehensive survey on this subject, please refer to Crane et al. (2020). Our work is specifically inspired from the ISOMAP method (Tenenbaum et al., 2000), a dimensionality reduction method which approximates geodesic

<sup>3</sup>This metric should be called *Log-probabilities*. However, since Log-odds is a commonly used name in the literature, we refer to it as Log-odds.

385 paths on a manifold. However, our method differs from  
 386 ISOMAP in that we weight our graph not using euclidean  
 387 distance, but the norms of a model’s gradients. Our aim is  
 388 indeed not to model the input space, but to explain a neural  
 389 network by building paths avoiding high-gradients regions.  
 390

391 As mentioned in Section 3, the idea of using a kNN algo-  
 392 rithm to avoid computing gradients on out of distribution  
 393 data points has also been used in Enhanced Integrated  
 394 Gradients Jha et al. (2020). However, this method creates a path  
 395 which is model agnostic, as it does not necessarily avoid  
 396 high gradients regions. As a result, it can lead to significant  
 397 artefacts which do not reflect the model’s behaviour. To  
 398 support this argument, we provide an example where this  
 399 method fails on the half-moons datasets (Figure 7).



400  
 401 **Figure 7. Enhanced IG attributions.** Enhanced IG computes a  
 402 kNN algorithm, uses Dijkstra to find the shortest path between an  
 403 input and a reference baseline, and computes gradients along this  
 404 path. However, this method is model agnostic and can as a result  
 405 cross a high gradients region, which is the case in this example,  
 406 between the input A and the baseline B. Input A therefore has a  
 407 high attribution which does not reflect the model’s true behavior.  
 408 In this example, the noise is  $\mathcal{N}(0, 0.15)$ .  
 409  
 410

411 The idea of adapting the path according to the model has been  
 412 proposed by Kapishnikov et al. (2021), calling their method  
 413 Guided Integrated Gradients. Their method computes this  
 414 path greedily by selecting around 10% of the features that  
 415 have the lowest absolute value of the partial derivatives, and  
 416 switching these features from the baseline’s values to the  
 417 input’s values. However, as the authors indicate, such method  
 418 can create out of distribution data points, with part of their  
 419 features either matching the input or the baseline. As a result,  
 420 they add a hyperparameter K, which forces the path to go  
 421 through K points along the straight line between the input  
 422 and the baseline. We argue, however, that, by directly approximating  
 423 the path of least resistance, our method is more principled compared  
 424 to Guided IG. In Guided IG, it is not clear how to choose the value of K: a too low value could cre-

425 ate out of distribution samples, while a too high one would  
 426 force the path to be close to the straight line and potentially  
 427 cross high gradients regions. It is similarly not clear why  
 428 switching 10% of the features, and how to tune this hyper-  
 429 paramer. On the other hand, we argue that Geodesic IG does  
 430 not have this issue. It has indeed two hyperparamers: the  
 431 number of points used to approximate the geodesic path, and  
 432 the number k of nearest neighbours. Yet, the performance  
 433 of Geodesic IG should improve when both hyperparamers  
 434 values increase, as, when more points or more neighbours  
 435 are used, the approximation of the geodesic path should  
 436 improve. Increasing these values would, however, require  
 437 more computing and performance needs to be balanced with  
 438 the amount of compute available.  
 439

## 5. Discussion

We have identified in this paper potential issues with path-based attribution methods: the presence of artefacts due to ignoring the curvatures of the model. To overcome these issues, we have introduced Geodesic Integrated Gradients, an adaptation of the original IG method which integrates gradients not along a straight line, but along the geodesic of a manifold defined by the model.

By avoiding high-gradients regions in the input space, we have shown that Geodesic IG can successfully address these issues. Moreover, it follows all of the axioms defined by Sundararajan et al. (2017).

We have presented two methods to approximate these geodesic paths: one based on kNN and the other Stochastic Variational Inference. Whilst our method shows clear advantage over the alternatives, evaluated by metrics such as Comprehensiveness and Log-Odds, there are some further research questions that we leave for future explorations. One is around computational intensiveness of our methods, as we discussed in section 3. The other is the inherent noise in any sampling process that we may wish to use for approximating the geodesic lines. Regarding the noise, it is conceivable that a method involving directly solving the geodesic equation might result in less noisy approximation of the geodesic paths. Depending on the way one solves such equations, it is possible that finding the geodesic paths might become more computationally efficient than using SVI.

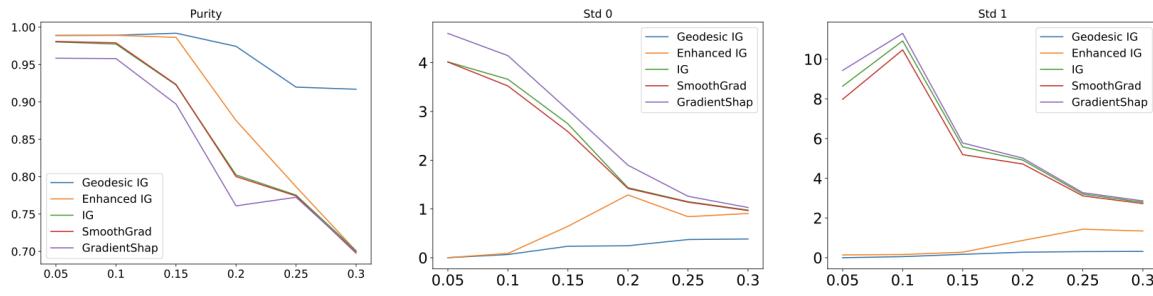
## References

- Chen, N., Ferroni, F., Klushyn, A., Paraschos, A., Bayer, J., and Smagt, P. v. d. Fast approximate geodesics for deep generative models. In *International Conference on Artificial Neural Networks*, pp. 554–566. Springer, 2019.
- Crane, K., Livesu, M., Puppo, E., and Qin, Y. A survey

- 440 of algorithms for geodesic paths and distances. *arXiv preprint arXiv:2007.10430*, 2020.
- 441
- 442 Das, A. and Rad, P. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371*, 2020.
- 443
- 444
- 445 DeYoung, J., Jain, S., Rajani, N. F., Lehman, E., Xiong, C., Socher, R., and Wallace, B. C. Eraser: A benchmark to evaluate rationalized nlp models. *arXiv preprint arXiv:1911.03429*, 2019.
- 446
- 447 Dombrowski, A.-K., Alber, M., Anders, C., Ackermann, M., Müller, K.-R., and Kessel, P. Explanations can be manipulated and geometry is to blame. *Advances in Neural Information Processing Systems*, 32, 2019.
- 448
- 449
- 450 Everingham, M., Van Gool, L., Williams, C. K. I., Winn, J., and Zisserman, A. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- 451
- 452 Hassija, V., Chamola, V., Mahapatra, A., Singal, A., Goel, D., Huang, K., Scardapane, S., Spinelli, I., Mahmud, M., and Hussain, A. Interpreting black-box models: a review on explainable artificial intelligence. *Cognitive Computation*, 16(1):45–74, 2024.
- 453
- 454
- 455 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- 456
- 457 Jha, A., K Aicher, J., R Gazzara, M., Singh, D., and Barash, Y. Enhanced integrated gradients: improving interpretability of deep learning models using splicing codes as a case study. *Genome biology*, 21(1):1–22, 2020.
- 458
- 459
- 460 Kapishnikov, A., Venugopalan, S., Avci, B., Wedin, B., Terry, M., and Bolukbasi, T. Guided integrated gradients: An adaptive path method for removing noise. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5050–5058, 2021.
- 461
- 462 Lundberg, S. M. and Lee, S.-I. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- 463
- 464
- 465 Marcus, G. Deep learning: A critical appraisal. *arXiv preprint arXiv:1801.00631*, 2018.
- 466
- 467 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- 468
- 469
- 470
- 471 Ribeiro, M. T., Singh, S., and Guestrin, C. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- 472
- 473 Sap, M., Card, D., Gabriel, S., Choi, Y., and Smith, N. A. The risk of racial bias in hate speech detection. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pp. 1668–1678, 2019.
- 474
- 475 Shrikumar, A., Greenside, P., Shcherbina, A., and Kundaje, A. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.
- 476
- 477 Shrikumar, A., Greenside, P., and Kundaje, A. Learning important features through propagating activation differences. In *International conference on machine learning*, pp. 3145–3153. PMLR, 2017.
- 478
- 479 Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- 480
- 481 Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. In *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- 482
- 483 Tenenbaum, J. B., Silva, V. d., and Langford, J. C. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- 484
- 485 Tonekaboni, S., Joshi, S., Campbell, K., Duvenaud, D. K., and Goldenberg, A. What went wrong and when? instance-wise feature importance for time-series black-box models. *Advances in Neural Information Processing Systems*, 33:799–809, 2020.
- 486
- 487 Yang, T., Arvanitidis, G., Fu, D., Li, X., and Hauberg, S. Geodesic clustering in deep generative models. *arXiv preprint arXiv:1809.04747*, 2018.
- 488
- 489 Zeiler, M. D. and Fergus, R. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pp. 818–833. Springer, 2014.
- 490
- 491
- 492
- 493
- 494
- 495
- 496
- 497
- 498
- 499

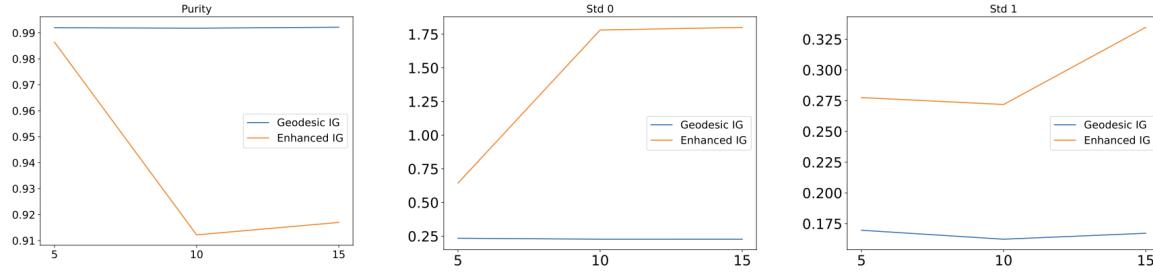
495 **A. Additional Half-moons results**

496 We present on Figure 8 more results on the half-moons dataset, using different amounts of noise. We can see that, for certain  
 497 amount of noise, Enhanced IG dramatically fails, while Geodesic IG performs consistently well on every amount of noise  
 498 tested here. We believe that the failure of Enhanced IG in the high-noise setting is due to the following reason. As noise  
 499 increases the points on either moons get closer to each other. As a result, the model loses the property that gradients are only  
 500 large on the decision boundary and fall rapidly as we move away. Therefore a lot of points are on the high-gradient regions.  
 501 However, Enhanced IG chooses the path purely based on nearest neighbours, ignoring model gradients. Hence, leading to  
 502 low purity. This is contrast to geodesic IG, which actively avoids regions of high gradients.



516 *Figure 8.* Evaluation of different attribution methods on the half-moons dataset with different amounts of noise:  $\mathcal{N}(0, x)$  where  $x$  is  
 517 defined as the axis of each plot.

518 We also perform here an ablation study using different values of  $k$  for the kNN algorithm. We show the results on Figure 9.  
 519 The results show that increasing  $k$  harms the performance of Enhanced IG, leaving the ones of Geodesic IG unchanged.  
 520 This is probably due to the fact that increasing  $k$  allows connections between points further apart, potentially crossing  
 521 high-gradients regions. While Geodesic IG would not follow such paths, Enhanced IG only uses euclidean distance, and is  
 522 therefore more likely to generate paths crossing high-gradients regions.



526 *Figure 9.* Evaluation of Geodesic IG and Enhanced IG for different values of  $k$  in the kNN algorithm. We can see that increasing this  
 527 parameter harms Enhanced IG performance, while it does not seem to have a major effect on Geodesic IG performance.

## B. Additional heatmaps and results on Pascal VOC 2012

We also qualitatively compare on Figure 10 Geodesic IG with the original IG on 5 different images of the Pascal VOC 2012 dataset. Geodesic IG heatmaps appears to be less blurry than the ones generated with original IG method.

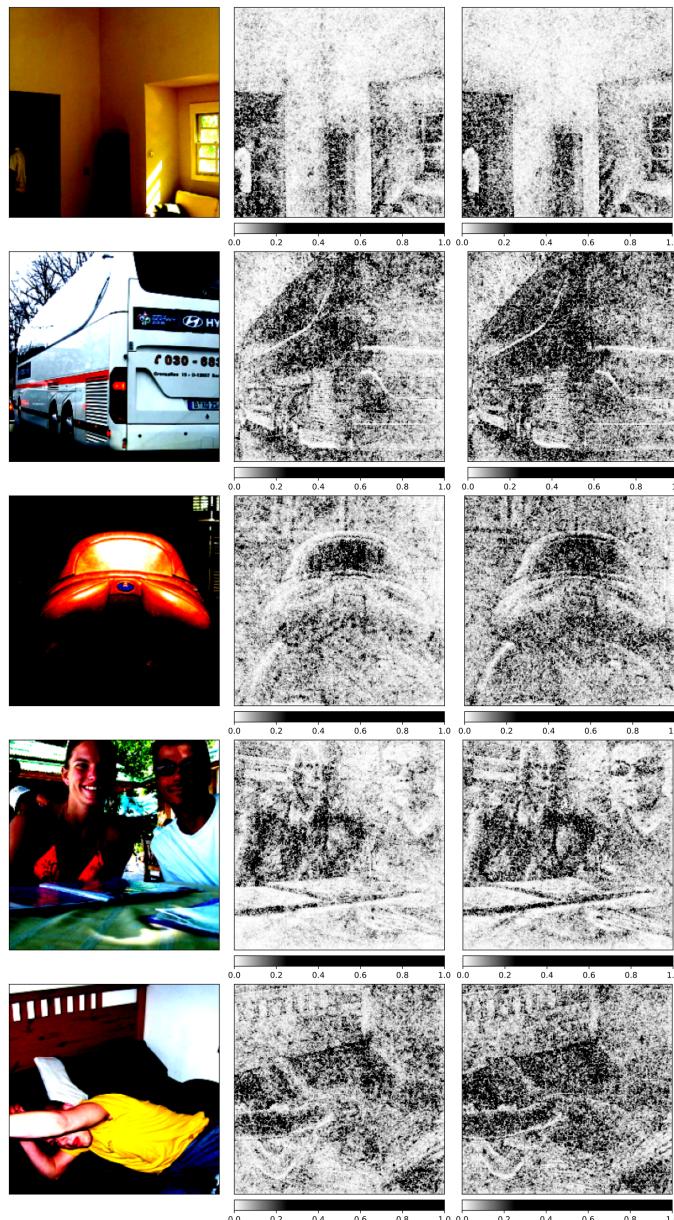


Figure 10. Heatmaps of Integrated Gradients (middle) and Geodesic IG (right) on 5 randomly chosen images from the test set of Pascal VOC 2012. We can see that Geodesic IG heatmaps are sharper compared with the original IG method.