
Harmonic Machine Learning Models are Robust

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce Harmonic Robustness, a powerful and intuitive method to test the
2 robustness of any machine-learning model either during training or in black-box
3 real-time inference monitoring without ground-truth labels. It is based on functional
4 deviation from the harmonic mean-value property, indicating instability and lack
5 of explainability. We show implementation examples in low-dimensional trees and
6 feedforward NNs, where the method reliably identifies overfitting, as well as in
7 more complex high-dimensional models such as ResNet-50 and Vision Transformer
8 where it efficiently measures adversarial vulnerability across image classes.

9 1 Motivation and Introduction

10 Modern application of Machine Learning (ML) across all industries faces numerous challenges in
11 maintaining quality of predictions: from the training phase where one must choose the “best” model
12 within a sea of architectures and hyperparameters to maximize performance without overfitting the
13 training data, while balancing with explainability and fairness in the context of Responsible AI; to
14 the inference phase where, in the face of production latency and throughput constraints, one must
15 efficiently monitor for performance degradation due to data drift; ideally this latter triggers the model
16 re-training phase, where one must revisit the model with freshly-labelled data, which, however, in
17 many applications such as credit card fraud may not be available till after a significant time lapse,
18 sometimes months later.

19 To address these challenges, we propose a simple geometric technique which enjoys several mitigating
20 properties:

- 21 • It is model-agnostic (black-box) and unsupervised, i.e., requiring no knowledge of model
22 inner workings, ground-truth labels or other auxiliary data.
- 23 • Its computation is algorithmically simple, linear in the number of data points tested, and has
24 good statistical sampling convergence.
- 25 • It reliably measures relative overfitting between two models on the same training data.
- 26 • It precisely measures model robustness across feature space and can immediately indicate
27 data drift in online monitoring.
- 28 • It is indicative of model explainability.

29 The particular proposal is to measure the “harmonicity” of the model, specifically the degree to
30 which the model function f satisfies the harmonic property,

$$\nabla^2 f = 0 \tag{1}$$

31 Functions which satisfy (1), i.e. “harmonic functions”, occur frequently in physics as solutions to
32 equilibrium problems involving minimization of energy, e.g. soap bubbles stretched on a boundary,
33 electrostatic field configurations, and heat flow (see Fig 1). They form smooth, minimal interpolations

34 between boundary values, and most importantly to our present discussion exhibit the “mean-value
35 property”,

$$f(x) = \frac{1}{S} \int_{B(x)} f d\Omega \quad (2)$$

36 which, in plain English, says that the value of the function at any point is the surface average of
37 the function over a ball of any radius surrounding the point (incidentally, (1) and (2) are equivalent
38 definitions)¹. The metric we propose, “anharmonicity” or γ for brevity, measures how well (2) is
39 satisfied over feature space, computing the difference between the function and its ball-averaged
40 value:

$$\gamma(x) \equiv |f(x) - \frac{1}{S_{r,n}} \int_{B(x,r)} f d\Omega_{r,n}| \quad (3)$$

41 where now we explicitly introduce the parameters r (ball radius) and n (feature dimension) as the
42 implementation details will depend on these. As the behavior of $f(x)$ may vary wildly over feature
43 space, so will $\gamma(x)$ depending on the degree to which f behaves in accordance with (2) for some
44 reasonable fixed choice of r . By association this will indicate which regions of feature space are
45 “more harmonic” for this model, and the average value of $\gamma(x)$ over feature space provides a summary
46 “anharmonicity” metric.

47 At this point, we can verify the above claimed mitigating properties of this metric:

- 48 • Measuring γ as per (3) requires nothing more than black-box access to the model, as is
49 expedient in an inference setting; this also allows testing of closed-source models without
50 having to request access to model details, facilitating efficiency of testing and helping to
51 keep the industry honest.
- 52 • As γ only requires computing the average value of f at a number of data points
53 approximating a ball, this is linear in the number of points and amenable to sampling.
- 54 • γ is proportional to the complexity of the decision surface; in particular for a binary classifier
55 it is proportional to the length of the decision boundary, which is positively correlated with
56 overfitting (see Appendix).
- 57 • If γ changes over time in online events, there must be data drift; this is great for online
58 monitoring where it is paramount to raise alerts as soon as a performance issue arises. If
59 production monitoring shows an increase in γ , one can pinpoint the data points responsible
60 and investigate that region of feature space more fully for counterfactuals — this might
61 trigger the need for re-training with more data or modeling in that region.
- 62 • Harmonic functions are natively explainable since, by the mean-value property, the
63 ‘explanation’ of any point is that it is the average of the points around it, which in turn are
64 ‘explained’ by their neighboring points, etc., all the way up to the feature boundaries which
65 have values fixed by some standard. The premier example is the linear function, of course
66 trivially harmonic by (1) and explainable by direct proportionality. Thus, the closer γ is
67 to zero, the more explainable the model will be. Conversely, the more a model fails (2) at
68 some point the more difficult it may be to explain, e.g., in the fraud domain if the average of
69 several non-fraudulent events was predicted to be fraudulent.

70 In short, γ is a proxy for the measure of model robustness in stability of prediction, resilience to data
71 drift, and ease of explainability.

72 Note, however, that the aim of this programme is certainly not to have a model be completely
73 harmonic — indeed by (2) it is easy to check that pure harmonic functions can have no local minima
74 or maxima, and that is too restrictive for real-world models.

75 Yet, we believe a robust, explainable ML model should be at least locally close to harmonic in
76 *most* of feature space², especially for production applications where stability is business-critical.

¹An easy informal way to see this equivalence is to recognize (1) as the divergence of the gradient ($\nabla^2 f = \nabla \cdot \nabla f$): the gradient expressing the change of f in all directions, if its divergence is zero then ‘change in f ’ neither flows into or out of any given point, hence the average change of f on any ball around that point is 0, relative to its value at the point.

²This can be made precise by choosing a value of r that is ‘small’ relative to typical distance between data points. Though as we’ll see below, the exact choice is not critical.

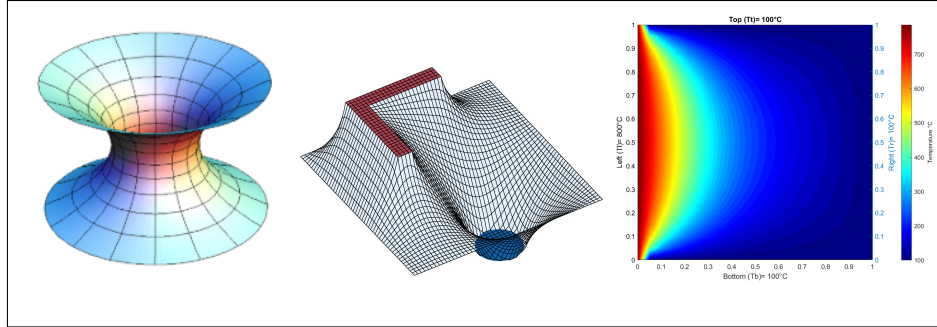


Figure 1: Examples of harmonic functions that appear in Nature: soap films[1], electrostatic potentials[2], and heat flows[3].

77 Harmonic functions exhibit the minimal curvature necessary to interpolate between fixed data points
 78 (see Appendix), and hence satisfy a certain Occam’s Razor of machine learning. Most real-world
 79 models will of course deviate from pure harmonicity, but γ gives us a way to track and quantify the
 80 deviation.

81 Finally, as far as we know, this is the first instance in the literature of an overarching, limiting
 82 *algebraic* standard on a ML function for the purpose of quality and stability. Our choice of the
 83 harmonic property is not in any way sacrosanct, but it is intuitive and accessible for quick and direct
 84 testing via the mean-value property (2), giving correlation with stability and ease of interpretation. It
 85 may be that some other class of functional algebraic constraints also captures this and more, a topic
 86 we leave open to the community to explore.

87 2 Related Work

88 Techniques to measure goodness of a predictive model are of course as old as the field of Machine
 89 Learning itself, traditionally centered on time-tested metrics such as precision, recall, F-score, AUC,
 90 etc. where the ground-truth labels of a test set are known. On the other hand, we are chiefly concerned
 91 with the real-world problem of measuring model robustness without access to ground truth labels,
 92 as for example occurs in a purely online inference environment with only black-box access to the
 93 model in conjunction with a live data stream. For this is the real, minimal environment in which most
 94 practitioners and end users of ML operate. Statistical techniques such as outlier or anomaly-detection
 95 [4] and distributional shift [5] enjoy usage here to give important hints of, but not true indications of,
 96 model robustness. Gradient-based methods such as PDP and ICE [6] where one looks for sudden
 97 changes in the decision function over feature space likewise may provide hints of robustness changes,
 98 though this differs from the current proposal as γ is measuring more than just the local feature
 99 sensitivity in the function, which may in fact be proper and desired behavior for, e.g., a steep linear
 100 response; rather, γ is measuring departure from *explainable* sensitivity as one sees in harmonic
 101 functions obeying the mean-value property (2).

102 Existing work [7] as well as a recent survey [8] reviews 23 metrics which are useful in the online
 103 inference setting, e.g., Average Confidence of True Class (ACTC) and Noise Tolerance Estimation
 104 (NTE). Within the black-box setting the metrics are essentially measuring how readily the predicted
 105 class label changes across feature space either due to targeted gradient-based search or random
 106 perturbation. These fit in the realm of Adversarial Machine Learning [9] which has blossomed into
 107 its own subfield, quite rightly dedicated to understanding the vulnerability of popular ML models
 108 to attacks [10][11] based on perturbing input data points. Recent work has found that adversarial
 109 weakness becomes more prevalent with increasing number of feature dimensions [12][13], and
 110 sensitive data domain dependence arises; unsurprisingly this is most apparent in image classification
 111 tasks [14], the premier testing-ground of adversarial ML, as each data point can easily contain
 112 thousands to millions of features (pixels); it is important to note, however, that many other domains,
 113 e.g., financial modeling, can contain thousands or more features and likewise be highly vulnerable
 114 to attacks [15]. Adversarial analyses also typically focus on class label changes, differing from our
 115 metric which is more precisely measuring the numerical stability of the prediction (logit) according

116 to the standard of harmonic geometry, and not merely the crossing of such predictions over discrete
117 thresholds leading to class label changes.

118 There have for a number of years been works focusing on the relationship between stability and
119 geometry of the classifier [16][17][18]: what these works find is that there is a correlation between
120 adversarial weakness and curvature of the decision surface. This has even fueled investigation into a
121 new way of classification using the average or majority-vote in a hypercube neighborhood of each
122 point [19]. This is corroborated by the present study as well, for harmonic functions describe minimal
123 surfaces with constant mean curvature [20], hence should have minimal adversarial weakness.

124 3 Method

125 Computing $\gamma(x)$ for a model as per its mathematical definition (3) is an extremely simple and
126 straightforward procedure, which we detail below in pseudocode:

Algorithm 1 Computation of γ at a point x in feature space

```
1: procedure  $\gamma(x)$ 
2:   ballPoints  $\leftarrow$  Ball( $x, r$ )
3:   N  $\leftarrow$  size(ballPoints)
4:   ballValue  $\leftarrow$  0
5:   for each point in ballPoints do
6:     ballValue  $+=$   $f(\text{point})$ 
7:   end for
8:   ballAvg  $\leftarrow$  ballValue/N
9:   return  $|f(x) - \text{ballAvg}|$ 
10: end procedure
```

127 One can then average $\gamma(x)$ over a region of feature space to get γ for that region, for example the
128 convex hull of a training set or all feature points seen in some inference production window.

129 In the above algorithm, the primary consideration is how to get the ball of radius r around the point x ,
130 i.e., Ball(x, r), remembering that in general x is a vector in some possibly high number of dimensions.
131 For any digital computation we will of course have to approximate a continuous ball with a discrete
132 number of points. The easiest-to-code solution is to construct a large number N of random vectors
133 around x , each normalized to some small magnitude r , hoping for isotropy and centrality (zero overall
134 vector sum). The random-walk behavior of N random vectors will however doom one to a \sqrt{N} bias
135 in one direction or another (see Appendix).

136 A better solution from a theoretical perspective is to form the "n-simplex" around each point. In two
137 dimensions, for example, the 2-simplex is an equilateral triangle; in three dimensions a tetrahedron,
138 etc. (see Appendix). In any arbitrary number of dimensions, the n -simplex centered about a point
139 will be maximally symmetric, hence ideally space-covering. One can further add balanced rotations
140 of the basic n -simplex, the more of which you add the closer the discrete approximation converges to
141 the continuous ball.

142 With a bit of linear algebra to compute the n -simplices (see Appendix), numerical trials indeed show
143 that n -simplices symmetrically cover space much more effectively than random vectors and, when
144 used to approximate the ball in our algorithm, accurately identify pure harmonic versus non-harmonic
145 functions, the details of which the interested reader may refer to in the Appendix. What we would
146 like to focus on in the remainder of this paper is actual ML models, as readers will find this most
147 applicable to their work.

148 4 Application to low-dimension models

149 For clarity in demonstrating the method let us first focus on basic models targeting a small, well-
150 understood dataset in the ML community: the Wine dataset [21], describing 13 features of three
151 different wines grown in the same region of Italy. We further restrict this analysis to just two of those
152 dimensions, "flavanoids" and "OD280/OD315 of diluted wines", in order to show a model taking a
153 two-dimensional input vector to a scalar output (the wine class label). What we seek to show is how

154 computing γ on models trained with these features on this data indicate the property of being well-fit
 155 or overfit, hence vulnerable to adversarial attacks or under-performance in production, purely from
 156 inference on out-of-training points without referencing any ground-truth labels.

157 We split the original data 80/20 to a Train/Test set, foregoing the usual split to a Validation set not
 158 just due to data sparsity (there only being around 100 data points in this set), but to show how the
 159 present technique can by itself indicate overfitting. We train four models: two Gradient Boosted
 160 Decision Trees (GBDT), "GBDT-1", with hyperparameters optimized on a grid search with 10-fold
 161 cross-validation, and the second, "GBDT-2", chosen with more extreme hyperparameter values; and
 162 two feedforward neural nets (multi-layer perceptron), a 1-hidden-layer model "MLP-1" and a much
 163 more parameterized 3-layer model "MLP-2" likely to overfit. Hyperparameters for all models are
 164 shown in Tables 1 and 2.

165 Intuition should tell us that the over-parameterized and under-regularized models GBDT-2 and MLP-2
 166 will perform better than GBDT-1 and MLP-1 on the Train set but not so on the Test set, and this is
 167 indeed the case as shown in the tables below. What we will show is this could also have been gleaned
 168 from the shapes of the decision boundaries for these classifiers, shown in the top row of Fig. 2: note
 169 how the shape of the overfit models' decision boundaries is more complicated than that of the well-fit
 170 models. Computing $\gamma(x)$ with $r = 0.05$ on a grid³ in a region safely enclosing all data points, we
 171 get the bottom row of plots. Notice that γ is non-zero only around the decision boundary, which for
 172 the overfit functions is always longer — these latter will thus have higher average γ , as we confirm
 173 in the tables. Computing γ thus identifies potential overfitting without need for checking a labeled
 174 validation or test set.

175 What's also interesting is that both wellfit models (GBDT-1 and MLP-1) have the same Test
 176 performance of 83%, but GBDT-1 has a slightly better γ (0.014) versus that of MLP-1 (0.016).
 177 One can easily see this difference from the lengths of the decision boundaries: the GBDT boundary
 178 consists of nearly straight lines while the MLP boundary is more curved. This illustrates how a model
 179 trainer with these test results might, from the perspective of robustness, prefer the GBDT model for
 180 use in inference.

Table 1: GBDT models

	GBDT-1	GBDT-2
max_depth	1	100
n_estimators	5	200
min_samples_split	2	2
learning_rate	0.1	1
Train accuracy	85%	100%
Test accuracy	83%	80%
γ (r=0.05)	0.014(2)	0.051(2)

Table 2: MLP models

	MLP-1	MLP-2
max_iter	200	1000
layer dims	(2,100,1)	(2,100,500,1000,1)
initial learning_rate	0.001	0.01
alpha	$1 \cdot 10^{-4}$	0
Train accuracy	82%	86%
Test accuracy	83%	79%
γ (r=0.05)	0.016(1)	0.027(1)

181 It should be clear from the above, then, that the Harmonic Robustness metric clearly works on simple
 182 ML functions in low numbers of dimensions, where we can visually confirm areas of feature space
 183 which are more robust and compare robustness of different models over the same space.

184 5 Application to high-dimension models

185 In this section we will consider the more complex case of higher-dimensional models to illustrate
 186 how the technique adapts.

187 The challenge of models over a larger number of dimensions (into the thousands or even millions) is
 188 three-fold:

- 189 1. High-dimensional simplices are more expensive in compute and storage.
- 190 2. High-dimensional models are typically more complex and take longer to run.
- 191 3. γ itself might be high-dimensional and interpretation is not straightforward.

³Dependence on r or grid-size was very mild; see Appendix.

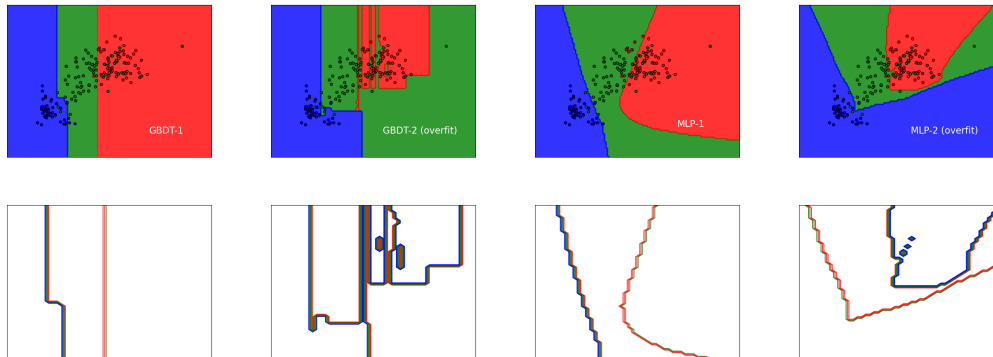


Figure 2: Decision regions (top row) and gamma contours (bottom row) of the two classifiers: GBDT (left), and MLP (right). Gamma contours for radius=0.05 closely follow decision boundaries. Boxed region shown corresponds to $[[0,5],[1,4]]$ in the x (“flavanoids”) and y (“OD280/OD315”) plane.

192 For the first challenge, Mathematics is actually kind to us, where it turns out that for large n the
 193 n -simplex is approximately the same as the vertices of the n -dimensional hypercube (see Appendix),
 194 and that is trivial to compute. For the second challenge, we will have to limit the number of points
 195 on the ball in order to be able to compute γ in a reasonable amount of time. Thus we may take a
 196 random sampling of the hypercube as a necessary approximation. Finally, if the output is not just
 197 a scalar, but rather multidimensional, one must decide whether some additional transformation is
 198 needed for interpretation. To take a weather example, if the model output is a 3-dimensional wind
 199 velocity, then $\vec{\gamma}$ represents the instability in velocity, and one might want to take its magnitude or
 200 angle with respect to north to interpret as speed instability or directional bias, respectively.

201 For the purpose of demonstration, we choose here to focus on high-dimensional image-classification
 202 models, due to popular practicality and ease of interpretation. The inputs (pixels) and outputs (class
 203 logits) are typically both high-dimensional, and would thus serve to illustrate the behavior of any
 204 other high-dimensional model as well. In particular, we consider ResNet-50 [22] and the Vision
 205 Transformer [23].

206 5.1 ResNet-50 and ViT

207 ResNet-50 and Vision Transformer (ViT) are image classifiers trained on 1000 distinct classes. To
 208 keep things manageable in this short work, we employ several restrictions: (1) Data is restricted to
 209 grayscale images: the value of every pixel is thus an integer from 0 to 255; (2) Images are rescaled to
 210 100x100 resolution: each image will thus be a 10000-dimensional vector; (3) Approximate simplices:
 211 10000-simplices are well-approximated as 1-hot vectors on the 10000-dim unit-hypercube as noted
 212 above; we scale each vector to magnitude 100 which amounts to a significant tone change at the
 213 position of the corresponding pixel. To increase ball coverage, we will use the simplices together
 214 with their reflections (anti-simplices); (4) Random sampling from simplices: rather than compute at
 215 all 20000 ball points (10000 simplex + 10000 anti-simplex) for each image, we take a random 0.1%
 216 sampling of such; (5) We compute γ only in the logit occupying the dimension of its predicted class
 217 label.

218 We apply these models to an animals test set [24] consisting of over 20k color images of animals in
 219 various resolutions from a pre-determined set of 10 classes (dog, horse, elephant, butterfly, chicken,
 220 cat, cow, pig, spider, squirrel). Then, as described above, we rescale each image to 100x100 pixels
 221 and convert to grayscale before computing γ in its predicted class logit dimension. We evaluate 100
 222 images per class, which will be sufficient to see the trends in robustness and justify the use of our
 223 approximations above.

224 For each image, we also execute an adversarial search process wherein we follow the gradient of γ
 225 for 25 iterations, recording the final image and its predicted class (see Algorithm 2 below). Each
 226 such image being only 25 pixels disparate from its original form, a change in class label is interpreted
 227 as “instability” in the original image, reminiscent of earlier gradient-based adversarial work [25].

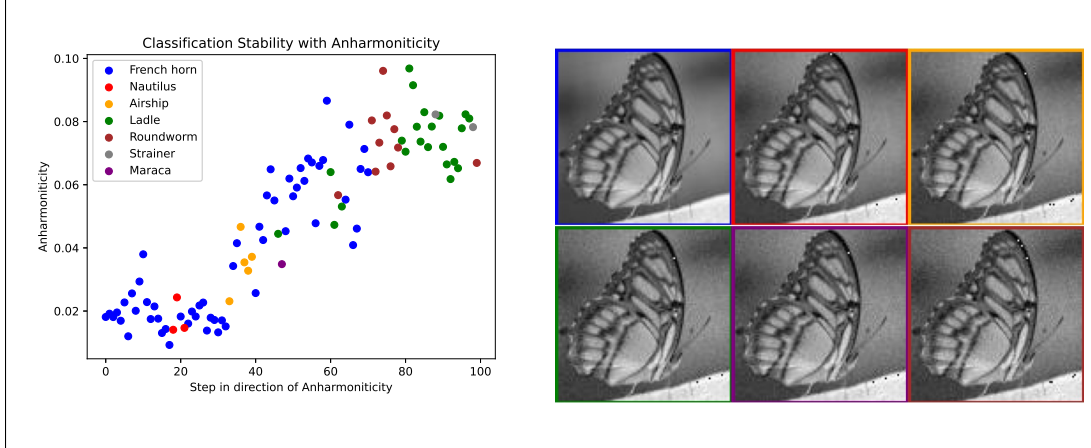


Figure 3: Demonstration of adversarial search procedure: following the stochastically increasing gradient of anharmonicity brings out classification instability. The image at step N differs from the original image by N pixels.

228 Note this procedure is actually stochastic gradient ascent as our γ -computation is based on random
 229 sampling of the hypercube.

Algorithm 2 γ -Stochastic Adversarial Search at a point x in feature space

```

1: procedure ADVERSARIALSEARCH( $x, r, N$ )
2:   currPoint  $\leftarrow x$ 
3:   numSteps  $\leftarrow N$ 
4:   for each step in numSteps do
5:     currGammas  $\leftarrow \{ \}$ 
6:     ballPoints  $\leftarrow \text{Ball}(\text{currPoint}, r)$ 
7:     for each point in ballPoints do
8:       currGammas[point]  $\leftarrow \gamma(\text{point}, r)$ 
9:     end for
10:    currPoint  $\leftarrow \text{argmax}(\text{currGammas})$ 
11:  end for
12:  return currPoint
13: end procedure

```

230 We chose 25 as the number of steps to execute as preliminary experiments showed this is generally
 231 the number of pixels one needs to change for the data and models under review before adversarial
 232 examples appear. Figure 3 shows one of such experiments where we execute the adversarial search
 233 for 100 steps, the original class label changing ever more frequently along that path of (stochastically)
 234 increasing γ . This procedure is actually very effective for quickly and reliably finding adversarial
 235 attacks on any input image, e.g., see Figure 4 where we show examples from each class where the
 236 predicted class radically changes after changing just 25 scattered pixels according to our adversarial
 237 scheme. Presumably these models behave more stably on larger color images, but it is useful to see
 238 how they behave on out-of-domain data, and γ gives you a way to measure that.

239 The reader may consult the Appendix for results on the full 1000 images, which we summarize
 240 here as supporting the conclusion that the generally more accurate ViT is also more robust(stable)
 241 than ResNet-50, except for the Cow and Squirrel class, and that γ , or rather γ in combination with
 242 the predicted class probability itself, accurately predicts this pattern as follows: the class softmax
 243 probability \mathcal{P}_C , controlled by the class logit (L_C) and average logit (\bar{L}),

$$\text{Prob}(\text{softmax})_C \equiv \mathcal{P}_C \equiv \frac{e^{L_C}}{\sum_{i=1..N_{classes}} e_i^L} \approx \frac{e^{L_C}}{e^{L_C} + (N_{classes} - 1) \cdot e^{\bar{L}}} \quad (4)$$

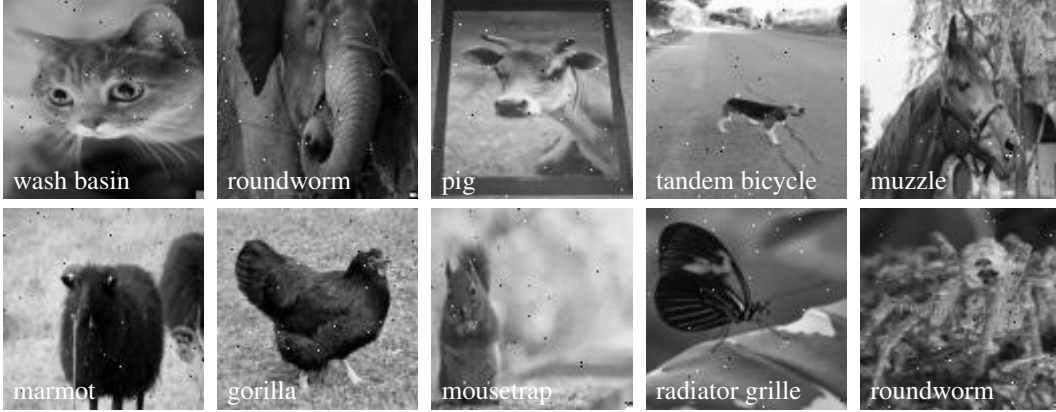


Figure 4: Examples of adversarial examples in ResNet-50 from following gradient of γ for 25 steps. Each image was originally correctly classified, but changed classes with modification of 25 scattered pixels as shown.

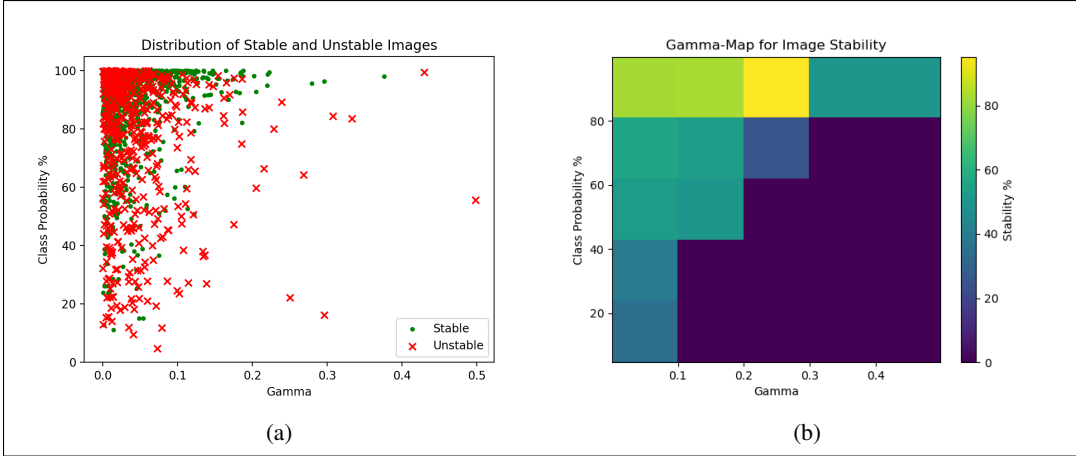


Figure 5: (a) Plotting predicted class probability \mathcal{P} and γ for 1000+1000 images classified with ResNet and ViT shows the unstable images tend to dominate high- γ /low- \mathcal{P} regions. (b) Density version of the previous plot serves as a practical "Gamma Map".

244 after a certain number N of gradient steps, is reduced by virtue of the class logit decreasing, on
 245 average, from L_C to $L_C - N\gamma$, so the adjusted probability becomes

$$\mathcal{P}'_C \approx \frac{e^{L_C - N\gamma}}{e^{L_C - N\gamma} + (N_{classes} - 1) \cdot e^{\bar{L}}} \approx \mathcal{P}_C e^{-N\gamma} \quad (5)$$

246 This metric, $\mathcal{P}_C e^{-N\gamma}$, gives a sort of "N-step adversarial robustness" which one can immediately
 247 measure at inference time and correlates well with actual image stability. We obtain visual
 248 confirmation of this metric by plotting the measured values of Class Probability \mathcal{P} and γ for each
 249 image, as well as whether it is stable or not after $N=25$ iterations, obtaining a "Gamma Map"
 250 (Figure 5). As a practical tool, this type of plot allows one to immediately gauge whether a predicted
 251 classification is likely to be stable just from measuring \mathcal{P} and γ , i.e., without having to do a full
 252 adversarial search.

253 6 Discussion and Conclusion

254 The foregoing demonstrated computation and interpretation of anharmonicity (γ) as a
 255 robustness(stability) metric for both a low-dimensional training setting, where γ can act as a regulator
 256 and rank models by degree-of-overfitting, as well as a rather different high-dimensional inference

257 setting where γ can feed into real-time performance feedback — from this the reader can interpolate
258 and extend to other scenarios.

259 Algorithmic implementation of γ is quite simple, straightforward, and applicable to any model
260 function as a measure of robustness. Monitoring systems and testing procedures can easily integrate
261 computations of γ as an alerting and regression test mechanism, respectively; for as we saw above for
262 ViT and ResNet, accuracy does not imply robustness. We see no reason why model builders should
263 hold back from computing γ alongside usual validation loss and other metrics to control overfitting.
264 As a proxy for explainability, this may lead into incorporating other metrics for Responsible AI into
265 the model life-cycle. As a metric to publish with a model’s quality card, one can envision reporting γ
266 for different data sets, indicating where a model is expected to more perform robustly. As a standard
267 for ML model quality, functional standards are translatable, shareable, and optimizable across the
268 industry; they may even point to certain mathematical truths pertaining to optimal ML systems.

269 We close with emphasizing a possibly trailblazing facet of our work: that one may apply a functional
270 mathematical standard, i.e., conformity to the properties of harmonic functions, to a ML system as
271 a way of assessing its quality and propriety for public usage. The harmonic standard may not be
272 necessarily ideal, most real-world ML functions being far from harmonic, but we posit that it is better
273 to reference a mathematically sound standard than having no such standard at all, giving AI systems
274 free reign in their inner complexity while relying on conventional external metrics like precision and
275 recall for quality control. For, assuming the white- or gray-box environment is not always going to be
276 available to us, if we do not devise multiple ways to check models’ inner complexity in a black-box
277 environment, we will be giving up too much control over what these systems may surprise us with.

278 References

- 279 [1] Frederik Brasz. Soap Films: Statics and Dynamics (available at
280 <https://www.princeton.edu/~stonelab/Teaching/FredBraszFinalPaper.pdf>). 2010.
- 281 [2] <https://www.britannica.com/science/electricity/Deriving-electric-field-from-potential>.
- 282 [3] <https://www.mathworks.com/matlabcentral/fileexchange/85073-2d-transient-heat-conduction>.
- 283 [4] Roel Bouman, Zaharah Bukhsh, and Tom Heskes. Unsupervised anomaly detection algorithms
284 on real-world data: how many do we need? *Journal of Machine Learning Research*, 25(105):1–
285 34, 2024.
- 286 [5] Olivia Wiles, Sven Goyal, Florian Stimberg, Sylvestre Alvisè-Rebuffi, Ira Ktena, Krishnamurthy
287 Dvijotham, and Taylan Cemgil. A fine-grained analysis on distribution shift. *arXiv preprint*
288 *arXiv:2110.11328*, 2021.
- 289 [6] Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil Pitkin. Peeking inside the black
290 box: Visualizing statistical learning with plots of individual conditional expectation. *journal of*
291 *Computational and Graphical Statistics*, 24(1):44–65, 2015.
- 292 [7] Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya Nori, and
293 Antonio Criminisi. Measuring neural net robustness with constraints. *Advances in neural*
294 *information processing systems*, 29, 2016.
- 295 [8] Jun Guo, Wei Bao, Jiakai Wang, Yuqing Ma, Xinghai Gao, Gang Xiao, Aishan Liu, Jian Dong,
296 Xianglong Liu, and Wenjun Wu. A comprehensive evaluation framework for deep model
297 robustness. *Pattern Recognition*, 137:109308, 2023.
- 298 [9] Joana C Costa, Tiago Roxo, Hugo Proença, and Pedro RM Inácio. How deep learning sees the
299 world: A survey on adversarial attacks & defenses. *arXiv preprint arXiv:2305.10862*, 2023.
- 300 [10] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian
301 Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint*
302 *arXiv:1312.6199*, 2013.
- 303 [11] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing
304 adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

- 305 [12] Beilun Wang, Ji Gao, and Yanjun Qi. A theoretical framework for robustness of (deep) classifiers
306 against adversarial examples. *arXiv preprint arXiv:1612.00334*, 2016.
- 307 [13] Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S Schoenholz, Maithra Raghu, Martin
308 Wattenberg, Ian Goodfellow, and G Brain. The relationship between high-dimensional geometry
309 and adversarial examples. *arXiv preprint arXiv:1801.02774*, 2018.
- 310 [14] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry.
311 Adversarially robust generalization requires more data. *Advances in neural information
312 processing systems*, 31, 2018.
- 313 [15] Micah Goldblum, Avi Schwarzschild, Ankit Patel, and Tom Goldstein. Adversarial attacks
314 on machine learning systems for high-frequency trading. In *Proceedings of the Second ACM
315 International Conference on AI in Finance*, pages 1–9, 2021.
- 316 [16] Constantine Caramanis, Shie Mannor, and Huan Xu. Robust optimization in machine learning.
317 2011.
- 318 [17] Alhussein Fawzi, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Robustness of
319 classifiers: from adversarial to random noise. *Advances in neural information processing
320 systems*, 29, 2016.
- 321 [18] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, Pascal Frossard, and Stefano
322 Soatto. Robustness of classifiers to universal perturbations: A geometric perspective. *arXiv
323 preprint arXiv:1705.09554*, 2017.
- 324 [19] Xiaoyu Cao and Neil Zhenqiang Gong. Mitigating evasion attacks to deep neural networks via
325 region-based classification. In *Proceedings of the 33rd Annual Computer Security Applications
326 Conference*, pages 278–287, 2017.
- 327 [20] Giorgio Talenti. A note on the gauss curvature of harmonic and minimal surfaces. *Pacific
328 Journal of Mathematics*, 101(2):477–492, 1982.
- 329 [21] Forina Mea. PARVUS—an extendible package for data exploration, Classification
330 and Correlation Institute of Pharmaceutical and Food Analysis and Technologies,
331 Via Brigata Salerno, 16147 Genoa, Italy. Data available at [https://scikit-
332 learn.org/stable/modules/generated/sklearn.datasets.load_wine.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_wine.html), 1991.
- 333 [22] <https://huggingface.co/microsoft/resnet-50>.
- 334 [23] https://huggingface.co/docs/transformers/model_doc/vit.
- 335 [24] <https://www.kaggle.com/datasets/alessiocorrado99/animals10>.
- 336 [25] Yuhang Wu, Sunpreet S Arora, Yanhong Wu, and Hao Yang. Beating attackers at their own
337 games: Adversarial example detection using adversarial gradient directions. In *Proceedings of
338 the AAAI Conference on Artificial Intelligence*, volume 35, pages 2969–2977, 2021.