

EVALUATION OF SIMILARITY-BASED EXPLANATIONS

Kamil Grudzień
Krystian Sztenderski





Plan

1. What are similarity based explanations
2. What similarity metrics are used
3. Evaluation criterias
4. Results



What are similarity based explanations



What are similarity based explanations

- Instance-based
- Similar instances as evidence to support model prediction
- “I(the model) think this image is cat because similar images I saw in the past were also cat”
- Very simple to understand without machine learning expertise
- Better than SHAP or LIME?

$$D = \{z_{train}^{(i)} = (x_{train}^{(i)}, y_{train}^{(i)})\}_{i=1}^N$$

$$z_{test} = (x_{test}, f(x_{test}))$$

$$\bar{z} = \operatorname{argmax}_{z_{train} \in D} R(z_{test}, z_{train})$$



**What similarity measures
are used**





What similarity metrics are used

- L2 Metric: $R_{\ell_2}(z, z') := -\|\phi(z) - \phi(z')\|^2$
- Cosine Metric: $R_{\cos}(z, z') := \cos(\phi(z), \phi(z'))$
- Dot Metric: $R_{\text{dot}}(z, z') := \langle \phi(z), \phi(z') \rangle$

ϕ can be any from:

$$\phi(z) = \mathbf{x} \qquad \phi(z) = \mathbf{h}^{\text{last}} \qquad \phi(z) = \mathbf{h}^{\text{all}}$$



Not only so simple- Gradient-based metrics

- **IF:** $R_{\text{IF}}(z, z') := \langle g_{\hat{\theta}}^z, H^{-1} g_{\hat{\theta}}^{z'} \rangle$
- **GD:** $R_{\text{GD}}(z, z') := \langle g_{\hat{\theta}}^z, g_{\hat{\theta}}^{z'} \rangle$
- **RIF:** $R_{\text{RIF}}(z, z') := \cos(H^{-\frac{1}{2}} g_{\hat{\theta}}^z, H^{-\frac{1}{2}} g_{\hat{\theta}}^{z'})$
- **GC:** $R_{\text{GC}}(z, z') := \cos(g_{\hat{\theta}}^z, g_{\hat{\theta}}^{z'})$
- **FK:** $R_{\text{FK}}(z, z') := \langle g_{\hat{\theta}}^z, I^{-1} g_{\hat{\theta}}^{z'} \rangle,$

Where H and I are the Hessian and Fisher information matrices of the loss $\mathcal{L}_{\text{train}}$



Evaluation criterias





Model Randomization Test

$$R(z_{test}, z_{train}^{(\pi_f(1))}) \geq R(z_{test}, z_{train}^{(\pi_f(2))}) \geq \dots \geq R(z_{test}, z_{train}^{(\pi_f(N))})$$

- R - relevance metric
- π_f - permutation of indices with decreasing relevance
- π_f and π_{f_rand} must have a small rank correlation
- Checks faithfulness

Identical Class Test

$$\arg \max_{\mathbf{z}=(\mathbf{x},y)\in\mathcal{D}} R(\mathbf{z}_{\text{test}}, \mathbf{z}) = (\bar{\mathbf{x}}, \bar{y}) \implies \bar{y} = \hat{y}_{\text{test}}$$



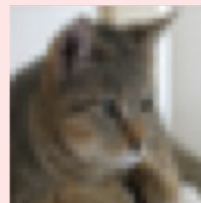
is cat because



a similar



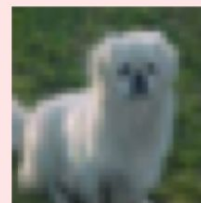
is cat.



is cat because



a similar

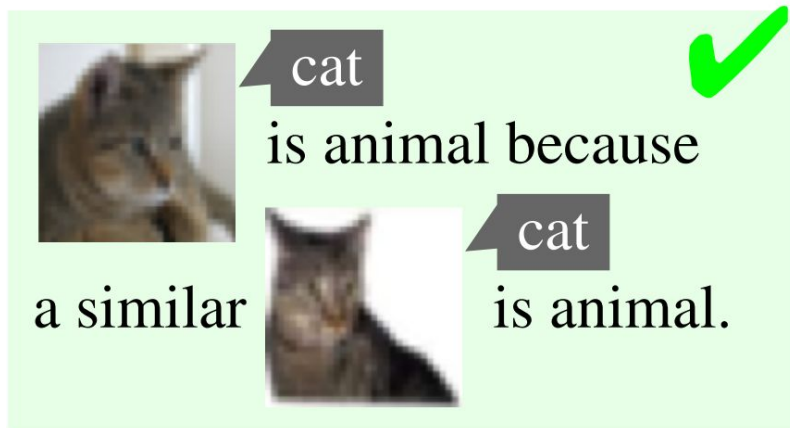


is dog.

Identical Subclass Test

$s(z)$ is a subclass for class y $s(z) \subset y$

$$\operatorname{argmax}_{z_{train} \in D} R(z_{test}, z_{train}) = \bar{z} \Rightarrow s(\bar{z}) = s(z_{test})$$

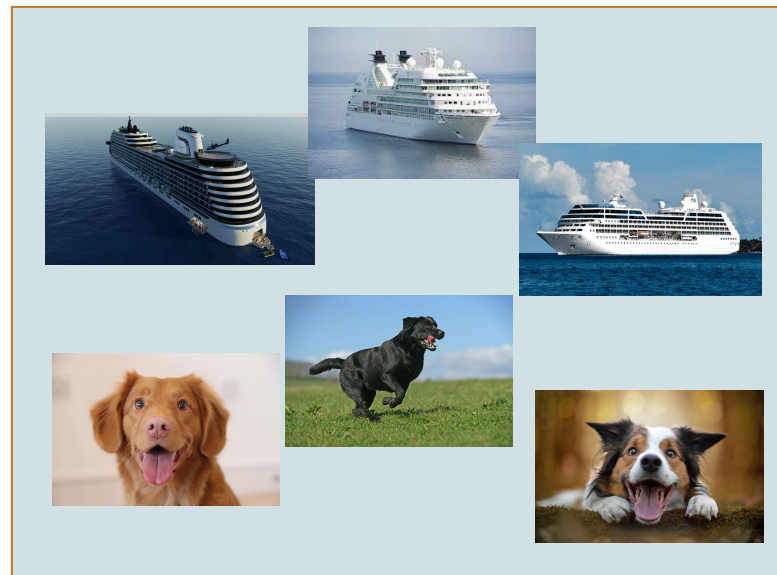


Dataset split

Class A



Class B

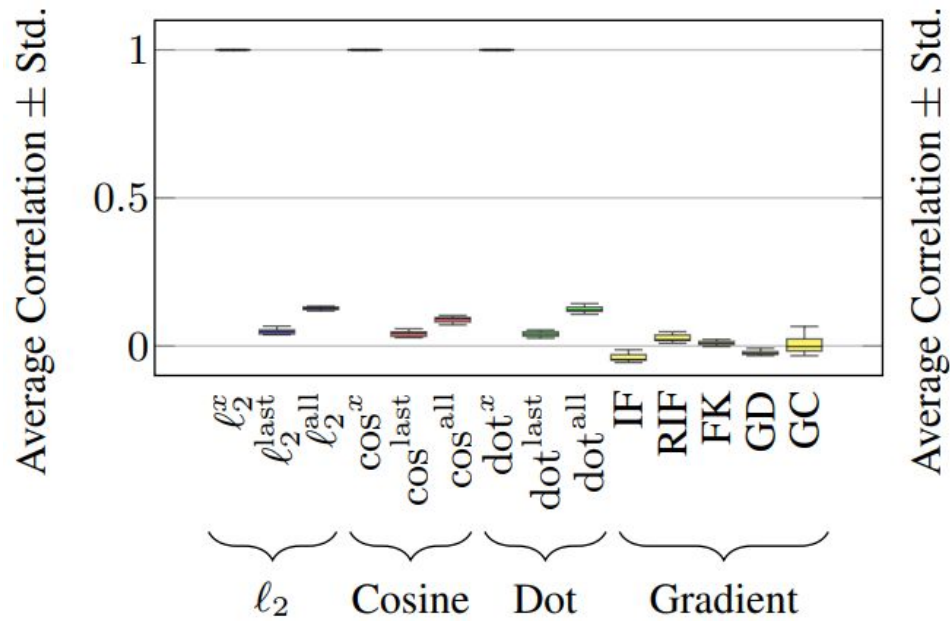




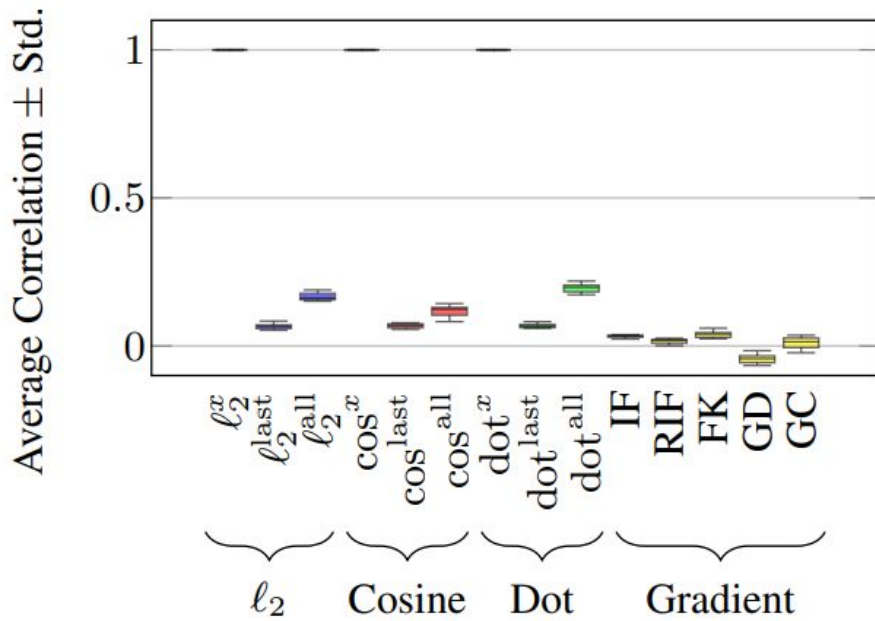
Results



Model randomization test

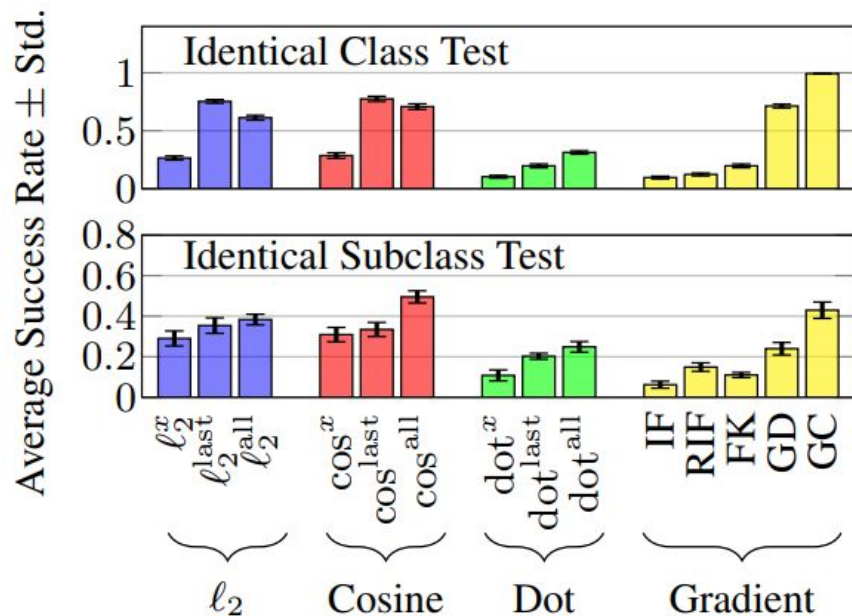


(a) CIFAR10 with CNN

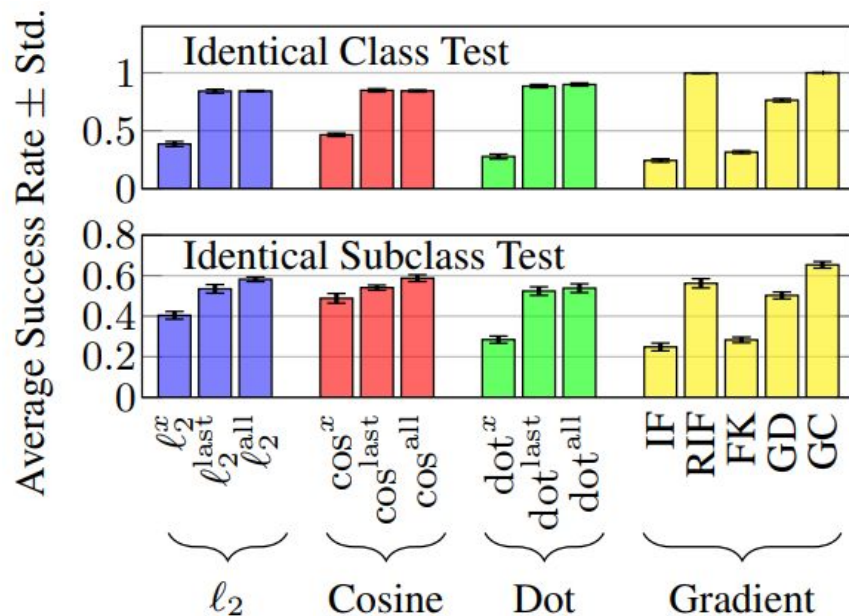


(b) AGNews with Bi-LSTM

Identical class and subclass test



















(a) CIFAR10 with CNN



(b) AGNews with Bi-LSTM

Frequently selected training instances

IF	Test Instances			Found
	airplane	frog	bird	truck
				
	$\text{COS}(\mathbf{z}_{\text{test}}, \mathbf{z}_{\text{train}})$			$\ \phi(\mathbf{z}_{\text{train}})\ $
	.00008	.00010	.00007	3,585
FK	Test Instances			Found
	cat	ship	bird	ship
				
	$\text{COS}(\mathbf{z}_{\text{test}}, \mathbf{z}_{\text{train}})$			$\ \phi(\mathbf{z}_{\text{train}})\ $
	.021	.020	.019	345,292,727
GD	Test Instances			Found
	cat	bird	horse	truck
				
	$\text{COS}(\mathbf{z}_{\text{test}}, \mathbf{z}_{\text{train}})$			$\ \phi(\mathbf{z}_{\text{train}})\ $
	.385	.291	.329	112.8
GC	Test Instances			Found
	truck	truck	truck	truck
				
	$\text{COS}(\mathbf{z}_{\text{test}}, \mathbf{z}_{\text{train}})$			$\ \phi(\mathbf{z}_{\text{train}})\ $
	.754	.754	.752	.0008

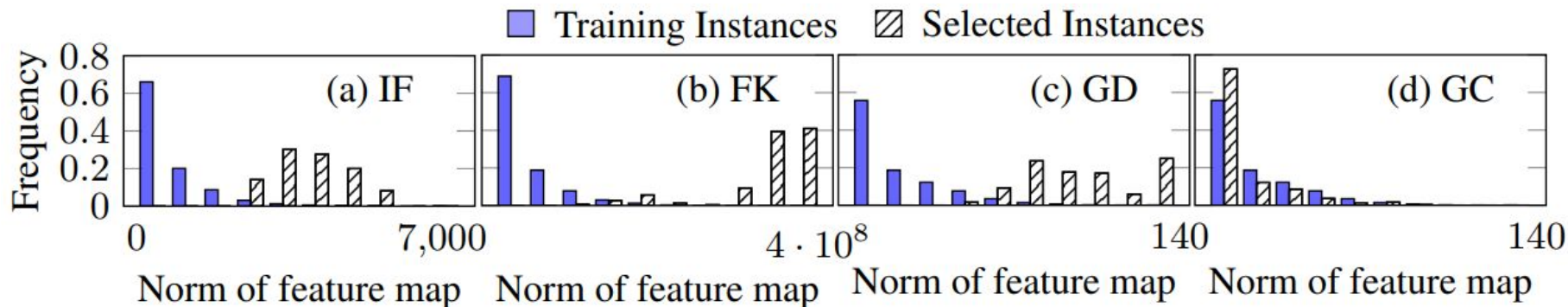


Failure of dot metrics & gradient based metrics

$$\langle \phi(\mathbf{z}_{\text{test}}), \phi(\mathbf{z}_{\text{train}}^{(i)}) \rangle < \langle \phi(\mathbf{z}_{\text{test}}), \phi(\mathbf{z}_{\text{train}}^{(j)}) \rangle$$

$$\|\phi(\mathbf{z}_{\text{train}}^{(i)})\| < \|\phi(\mathbf{z}_{\text{train}}^{(j)})\| \cos(\phi(\mathbf{z}_{\text{test}}), \phi(\mathbf{z}_{\text{train}}^{(j)}))$$

Frequently selected training instances





Takeaways

- 3 minimal requirement tests:

Model Randomization Test Identical Class Test Identical Subclass Test

- $\ell_2^{\text{last}}, \cos^{\text{last}}$ and **gradient-based** metrics **passed** the Model Randomization Test
- **Dot metrics** as well as **IF, FK and GD** **failed** the Identical Class and Subclass Test
- **GC** performed **the best** in most of the tests – **recommended** method



Thanks for your attention