## **Evaluation of Similarity-based Explanations**

Kazuaki Hanawa<sup>1,2</sup>, Sho Yokoi<sup>2,1</sup>, Satoshi Hara<sup>3</sup>, Kentaro Inui<sup>2,1</sup>

<sup>1</sup>RIKEN AIP, <sup>2</sup>Tohoku University, <sup>3</sup>Osaka University



## Analysis of the Failure

- Dot product-based methods do not work well
- Some instances are judged as similar to various test instances because of the large norm



Example of **Dot product of gradients**  $\langle g_{\text{test}}, g_i \rangle$ 



ship

cat



**g**<sub>test</sub>: Gradient of the test instance

horse

**g**<sub>i</sub>: Gradient of the *i*-th training instance

ship