

Learning to Explain with Complemental Examples

Atsushi Kanehira¹ and Tatsuya Harada^{2,3}

¹Preferred Networks, ²The University of Tokyo, ³RIKEN

Abstract

This paper addresses the generation of explanations with visual examples. Given an input sample, we build a system that not only classifies it to a specific category, but also outputs linguistic explanations and a set of visual examples that render the decision interpretable. Focusing especially on the complementarity of the multimodal information, i.e., linguistic and visual examples, we attempt to achieve it by maximizing the interaction information, which provides a natural definition of complementarity from an information theoretical viewpoint. We propose a novel framework to generate complemental explanations, on which the joint distribution of the variables to explain, and those to be explained is parameterized by three different neural networks: predictor, linguistic explainer, and example selector. Explanation models are trained collaboratively to maximize the interaction information to ensure the generated explanation are complemental to each other for the target. The results of experiments conducted on several datasets demonstrate the effectiveness of the proposed method.

1. Introduction

When we explain something to others, we often provide supporting examples. This is primarily because examples enable a concrete understanding of abstract explanations. With regard to machines, which are often required to justify their decision, do examples also help explanations?

This paper addresses the generation of visual explanations with visual examples. More specifically, given an input sample, we build a system that not only classifies it to a specific category but also outputs linguistic explanations and a set of examples that render the decision interpretable. An example output is shown in Fig. 1.

The first question to be raised toward this problem would be “How do examples help explanations?”, or equivalently, “Why are examples required for explanations?”

This work is done at the University of Tokyo.

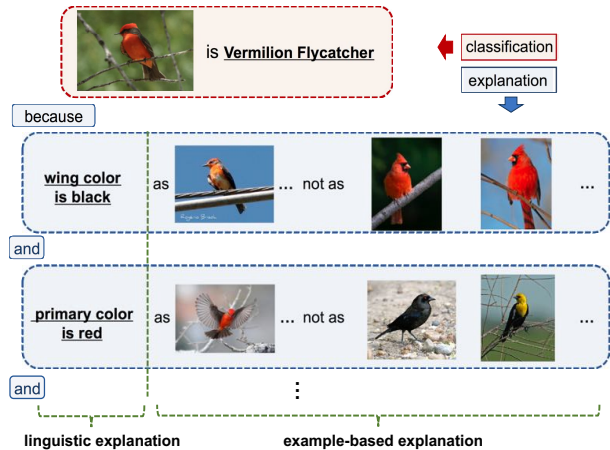


Figure 1: Our system not only classifies a given sample to a specific category (in the red dotted box), but also outputs linguistic explanations and a set of examples (in the blue dotted box).

To answer these questions, we consider the characteristics of two types of explanations pertaining to this work: linguistic explanation, and example-based explanation.

- Using language, one can transmit information efficiently by converting an event to a shared concept between humans. Inherently, the conversion process is invertible; thus, the whole event can not necessarily be represented by language alone.
- Using examples, one can transmit information more concretely than language can, as the saying, “a picture is worth a thousand words.” However, the way of the interpretation for the given examples is not determined uniquely. Thus, using examples alone is inappropriate for the explanation.

These explanations with different characteristics can be expected to complement each other, that is, from a lexicon, a thing that contributes extra features to something else in such a way as to improve or emphasize its quality [1].

The next important questions here are as follows: “How can the complementarity be achieved?” and “Which explanation is complemental, and which is not?”

We answer the former question from the information

theoretical viewpoint, that is, interaction-information [20] maximization. Interaction-information is one of the generalizations of mutual information defined on more than three random variables, and provides the natural definition of the complementarity: The increase of dependency of two variables when the third variable is conditioned.

We propose a novel framework in this work to build a system that generates complementary explanations. First, we introduce a linguistic explainer and an example selector parameterized by different neural networks, in addition to the predictor that is the target of the explanation. These two auxiliary models are responsible for generating explanations with linguistic and examples, respectively, and they are simultaneously trained to maximize the interaction information between variables of explanations and the output of the predictor in a post-hoc manner. Because the direct optimization of interaction-information with regard to the selector is intractable owing to the number of combination of examples, we maximize the variational lower bound instead. One more additional classifier, referred to as reasoner, appears in the computation of the lower bound. Taking linguistics and example-based explanations as inputs, the reasoner attempts to predict the output of the predictor. To enable the optimization of the selector with back-propagation, we utilized a reparameterization trick that replaces the sampling process of the examples with a differentiable function.

Under our framework, where complementarity is defined by information theory, we can understand better the complementary explanation related to the latter question. It can be mentioned that complementary examples for a linguistic explanation are a *discriminative set of examples*, by which one can reason to the correct conclusion with the given linguistic explanations, but cannot be achieved with different possible explanations. Complementary linguistic explanations to examples are also considered to be explanations that can construct such a set of examples. More details will be discussed in the subsequent section.

We conducted experiments on several datasets and demonstrated the effectiveness of the proposed method.

The contributions of this work are as follows:

- Propose a novel visual explanation task using linguistic and set of examples,
- Propose a novel framework for achieving complementarity on multimodal explanations.
- Demonstrate the effectiveness of the proposed method by quantitative and qualitative experiments.

The remainder of this paper is organized as follows. In Section 2, we discuss the related work of the visual explanation task. Further, we explain the proposed framework to achieve complementary explanations in Section 3 and describe and discuss the experiments that we performed on it in Section 4. Finally, we conclude our paper in Section 5.

2. Related Work

The visual cognitive ability of a machine has improved significantly primarily because of the recent development in deep-learning techniques. Owing to its high complexity, the decision process is inherently a black-box; therefore, many researchers have attempted to make a machine explain the reason for the decision to verify its trustability.

The primary stream is visualizing where the classifier weighs for its prediction by assigning an importance to each element in the input space, by propagating the prediction to the input space [24, 3, 29, 23, 30, 8, 31], or by learning the instance-wise importance of elements [4, 5, 15] with an auxiliary model. As a different stream, some works trained the generative model that outputs explanations with natural language [12, 21] in a post-hoc manner. Although most studies are focused on single modality, our work exploits multimodal information for explanations.

Prototype selection [25, 7, 13, 16, 25, 10] or machine teaching [18] can be considered as example-based explanations. The essential idea of these methods is to extract representative and discriminative (parts of) examples. In other words, they attempt to obtain examples that represent $p(\mathbf{x}|c)$, which is the distribution of sample \mathbf{x} conditioned on the category c . Our work is different in that we attempt to explain the black-box posterior distribution $p(c|\mathbf{x})$ such as that represented by deep CNN. Moreover, we utilized the linguistic information as well because the interpretation toward example-based explanation is not determined uniquely.

Few works have treated multimodality for explanation [21, 2], which is visual and linguistic. Although they provided visual information by referring to a part of the target samples, we explore the method to utilize other examples for explanation.

3. Method

The goal of this study is to build a model that generates linguistic and example-based explanations, which are complementary to each other. In this section, we describe the proposed framework. First, in subsection 3.1, we formulate our novel task with the notation used throughout this paper. Subsequently, the objective function to be optimized is illustrated in subsection 3.2. From subsection 3.3 to 3.6, we explain the details of the actual optimization process. The proposed method is discussed qualitatively in subsection 3.7. Finally, its relation with other explanation methods is mentioned in subsection 3.8.

3.1. Problem formulation

We denote by \mathbf{x} and \mathbf{y} the sample and the category that are the target for explanations where \mathbf{y} is a one-hot vector. \mathbf{s} is a vector representing a discrete attribute, whose

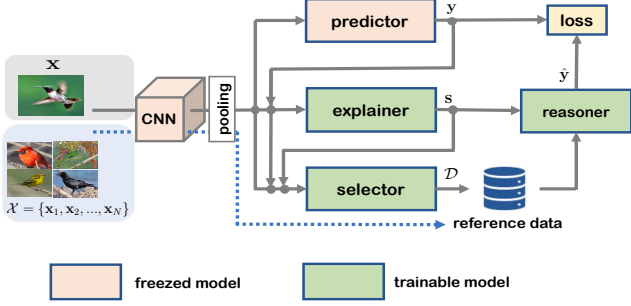


Figure 2: The pipeline of our explanation system. It holds two auxiliary models, which are responsible for generating explanations with linguistics and examples, respectively. In addition, it contains a reasoner that predicts the output of the predictor from the given explanations as described in subsection 3.3.

every index corresponds to the type of attribute (e.g., dim1 \rightarrow color, dim2 \rightarrow shape..), and the value of the vector corresponds to the value of attributes (e.g., 1 \rightarrow red, 2 \rightarrow blue ..). Attribute values are also treated as one-hot vector on implementation. We assume that the attributes are assigned to all the samples used for training the explanation model. In this study, we utilize one attribute as an element of linguistic explanation. More specifically, linguistic explanation s contains only one non-zero value (i.e., $\|s\|_0 = 1$), and the corresponding type-value is outputted (e.g., “because color is red.”). To explicitly distinguish the variable representing linguistic explanation from one representing attributes assigned to the samples, we denote the former by s and the latter by \hat{s} . The set of candidate examples used for explanation is represented by $\mathcal{X} = \{(x_i, \hat{s}_i, y_i)\}_{i=1}^N$, and its subset $\mathcal{D} \subset \mathcal{X}$, $|\mathcal{D}| = k$ is used as an element of the example-based explanation. We assume $\binom{N}{k}$, and that the number of combinations \mathcal{D} , is sufficiently large. Our system generates multiple elements $(s_1, \mathcal{D}_1), (s_2, \mathcal{D}_2), \dots, (s_M, \mathcal{D}_M)$, and construct a explanation by simply applying them to the template as in Fig. 1.

We built a model not only categorizing the input x to a specific class y , but also providing an explanation with linguistics and example-based explanations s and \mathcal{D} . We decomposed a joint distribution $p(y, s, \mathcal{D}|x)$ to three probabilistic models: predictor, explainer, selector, all of which were parameterized by different neural networks:

$$p(y, s, \mathcal{D}|x) = \underbrace{p(y|x)}_{\text{predictor}} \underbrace{p(s|x, y)}_{\text{explainer}} \underbrace{p(\mathcal{D}|x, y, s)}_{\text{selector}} \quad (1)$$

predictor $p(y|x)$ is the target model of the explanation, which categorizes sample x to y . Particularly, we study the model pretrained for the classification task. Throughout this paper, the weight of the predictor is frozen, and the remaining two auxiliary models, namely, explainer and selector, are trained to explain the output of the predictor.

explainer $p(s|x, y)$ is the probability of linguistic explanation s being selected given target sample x and class y . We limit $\|s\|_0 = 1$, and the dimension and the value corresponding to the non-zero element is used as an explanation.

selector $p(\mathcal{D}|x, y, s)$ is the probability of example-based explanation \mathcal{D} being selected out of all candidate examples given x, y , and s as inputs.

3.2. Objective function

We illustrate the objective function optimized for training the explanation models in this subsection. As stated earlier, linguistic explanation s and example-based explanation \mathcal{D} are expected to be complementary to each other. Intuitively, one type of explanation should contain the information for the target y , that is different from what the other explanation contains.

Hence, we leverage the interaction information [20] as an objective function. Interaction-information is a generalization of the mutual information defined on more than three random variables, and it measures how the dependency of two variable increases when the third variable is conditioned, which provides a natural definition toward complementarity.

From the definition, the interaction information of y, s, \mathcal{D} conditioned on the input x is written as the difference of two mutual information:

$$I(y, s, \mathcal{D}|x) = I(y, s|x, \mathcal{D}) - I(y, s|x) \quad (2)$$

where

$$\begin{aligned} & I(y, s|x, \mathcal{D}) \\ &= \int_{\mathbf{x}} \sum_{y, s, \mathcal{D}} p(y, s, \mathcal{D}, \mathbf{x}) \log \frac{p(y, s|x, \mathcal{D})}{p(s|x, \mathcal{D})p(y|x, \mathcal{D})} d\mathbf{x}, \\ &= \mathbb{E}_{p(\mathbf{x})} \left[\underbrace{\sum_{y, s, \mathcal{D}} p(s, \mathcal{D}|x, y) p(y|x) \log \frac{p(s|x, y, \mathcal{D})}{p(s|x, \mathcal{D})}}_{(A)} \right] \quad (3) \end{aligned}$$

and similarly,

$$\begin{aligned} & I(y, s|x) \\ &= \mathbb{E}_{p(\mathbf{x})} \left[\underbrace{\sum_{y, s} p(y, s|x) \log \frac{p(s|x, y)}{\sum_y p(s|x, y) p(y|x)}}_{(B)} \right] \quad (4) \end{aligned}$$

Intuitively, it measures how much linguistic explanation s becomes useful information to identify a category y when given a set of example-based explanation \mathcal{D} .

The direct estimation of (3) is difficult, as calculating the expectation over all possible \mathcal{D} is intractable. We handle

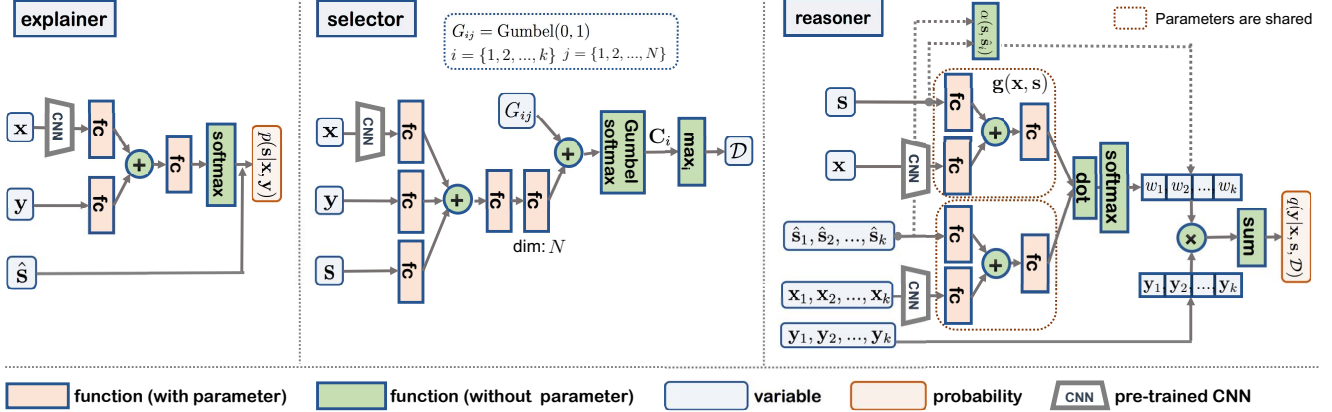


Figure 3: Structures of three neural networks representing three probabilistic models. As described in subsection 3.4, the network of the selector predict the parameter of categorical distribution unlike the other two models for the ease of optimization.

this problem by (a) introducing the variational lower bound and (b) leveraging reparameterization trick similar to [4], which are described in subsections 3.3 and 3.4, respectively.

3.3. Maximizing variational bound

In this subsection, we consider the variational lower bound of (A) in (3). From the definition of the KL divergence, $p \log p \geq p \log q$ is applied for any distribution p and q . Using this relation, (A) inside the expectation in (3) can be lower-bounded as follows:

$$(A) \geq \sum_{\mathbf{y}, \mathbf{s}, \mathcal{D}} p(\mathbf{s}, \mathcal{D} | \mathbf{x}, \mathbf{y}) p(\mathbf{y} | \mathbf{x}) \log \frac{q(\mathbf{s} | \mathbf{x}, \mathbf{y}, \mathcal{D})}{p(\mathbf{s} | \mathbf{x}, \mathcal{D})} \quad (5)$$

$q(\mathbf{s} | \mathbf{x}, \mathbf{y}, \mathcal{D})$ can be any distribution provided that it is normalized, and the expectation of the KL divergence between $q(\mathbf{s} | \mathbf{x}, \mathbf{y}, \mathcal{D})$ and true distribution $p(\mathbf{s} | \mathbf{x}, \mathbf{y}, \mathcal{D})$ is the difference between (A) and the lower bound. Similar to the method in [9], we used the following

$$q(\mathbf{s} | \mathbf{x}, \mathbf{y}, \mathcal{D}) = \frac{q(\mathbf{s}, \mathbf{y} | \mathbf{x}, \mathcal{D})}{q(\mathbf{y} | \mathbf{x}, \mathcal{D})} = \frac{q(\mathbf{y} | \mathbf{x}, \mathbf{s}, \mathcal{D}) p(\mathbf{s} | \mathbf{x}, \mathcal{D})}{\sum_{\mathbf{s}'} q(\mathbf{y} | \mathbf{x}, \mathbf{s}', \mathcal{D}) p(\mathbf{s}' | \mathbf{x}, \mathcal{D})},$$

and substituted it to (5). When considering parameterizing as in (1), it is computationally difficult to calculate $p(\mathbf{s} | \mathbf{x}, \mathcal{D})$. Considering the sampling order, we approximate it to $p(\mathbf{s} | \mathbf{x}, \mathbf{y})$ instead for simplicity. The first term of the objective function used for optimization is as follows:

$$(5) \approx \mathbb{E}_{p(\mathbf{y}, \mathbf{s}, \mathcal{D} | \mathbf{x})} \left[\log \frac{q(\mathbf{y} | \mathbf{x}, \mathbf{s}, \mathcal{D})}{\mathbb{E}_{p(\mathbf{s} | \mathbf{x}, \mathbf{y})} [q(\mathbf{y} | \mathbf{x}, \mathbf{s}, \mathcal{D})]} \right]. \quad (6)$$

$q(\mathbf{y} | \mathbf{x}, \mathbf{s}, \mathcal{D})$ is hereafter referred to as a reasoner, which is expected to *reason* the category of input given a pair of explanation for it.

3.4. Continuous relaxation of subset sampling

The abovementioned (6) is required to be optimized stochastically with sampling to avoid calculating the summation over the enormous number of possible combinations

of \mathcal{D} . In this situation, the difficulty of optimization with regard to the network parameter still exists. As it involves the expectation over the distribution to be optimized, sampling process disables calculating the gradient of parameters, rendering it impossible to apply back-propagation.

We resort on the reparameterization trick to overcome this issue, which replaces the non-differential sampling process to the deterministic estimation of the distribution parameter, followed by adding random noise. In particular, the Gumbel-softmax [19, 14] function is utilized similar to [4], which approximates a random variable represented as a one-hot vector sampled from a categorical distribution to a vector using continuous values. Specifically, we estimate the parameter of categorical distribution $\mathbf{p} \in R^N$ satisfying $\sum_{i=1}^N p_i = 1$ by the network where $N = |\mathcal{X}|$ is the candidate set of examples. An N -dimensional vector \mathbf{C} , which is a continuous approximation of the categorical one-hot vector, is sampled by applying softmax to the estimated parameter after taking logarithm and adding a noise \mathbf{G} sampled from the Gumbel distribution as follows:

$$\mathbf{C}[i] = \frac{\exp\{(\log p_i + G_i)/\tau\}}{\sum_{j=1}^N \exp\{(\log p_j + G_j)/\tau\}} \quad (7)$$

where

$$G_i = -\log(-\log u_i), u_i \sim \text{Uniform}(0, 1), \quad (8)$$

and τ is the temperature of softmax controlling the hardness of the approximation to the discrete vector. To sample k -hot vector representing example-based explanation \mathcal{D} , concrete vector \mathbf{C} is independently sampled k times, and element-wise maximum is taken to $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_k$ to construct a vector corresponding to \mathcal{D} .

3.5. Structure of networks

We parameterize three probabilistic distributions, explainer, selector, and reasoner with different neural networks. We elucidate their detailed structures.

Explainer $p(\mathbf{s}|\mathbf{x}, \mathbf{y})$ is represented by a neural network that predicts the probability of each type (dimension) of attribute is selected. The model is constituted using three fully-connected layers as in left of Fig. 2. Taking the target sample \mathbf{x} and the category label \mathbf{y} as inputs, the model projects them to the common-space and element-wise summation is applied. After one more projection, they are normalized by the softmax function. The output dimension of the network $\mathbf{f}(\mathbf{x}, \mathbf{y})$ is the same as that of the attribute vector, and each dimension indicates the probability that each type of attribute is selected as an explanation. When training, the attribute $\hat{\mathbf{s}}$ assigned to the sample is used as the value. Formally, for all i -th dimension of the linguistic explanation vector,

$$p(\mathbf{s}|\mathbf{x}, \mathbf{y}) = \begin{cases} \mathbf{f}(\mathbf{x}, \mathbf{y})[i] & \text{if } \mathbf{s}[i] = \hat{\mathbf{s}}[i] \\ 0 & \text{otherwise} \end{cases}$$

For inference, the value that maximizes the output of the reasoner (described later) for the class to be explained is selected.

Selector $p(\mathcal{D}|\mathbf{x}, \mathbf{y}, \mathbf{s})$ takes the linguistic explanation \mathbf{s} in addition to \mathbf{x} and \mathbf{y} as inputs; their element-wise summation is calculated after projecting them to the common-space. As stated in the previous subsection, we leverage reparameterization trick to render the optimization tractable owing to the enormous number of the combination \mathcal{D} . The network estimates the parameter \mathbf{p} of categorical distribution. When sampling from a distribution, noise variables that are independently generated k times are added to the parameter, and the element-wise maximum is computed after the Gumbel softmax is applied.

Reasoner $q(\mathbf{y}|\mathbf{x}, \mathbf{s}, \mathcal{D})$ infers the category to which the sample \mathbf{x} belongs, given a pair of generated explanation $(\mathbf{s}, \mathcal{D})$. We design it by modifying the matching network [26], which is a standard example-based classification model. The prediction of the reasoner must be based on the given explanations. Such reasoning process is realized by considering (1) consistency to the linguistic explanation \mathbf{s} , and (2) similarity to the target sample \mathbf{x} , for each example in \mathcal{D} .

Based on a certain reason, the reasoner decides whether each example deserves consideration, and predicts the category exploiting only selected examples. The weight of each referred sample \mathbf{x}_i is determined by the visual and semantic similarity to the target \mathbf{x} . More formally,

$$q(\mathbf{y}|\mathbf{x}, \mathbf{s}, \mathcal{D}) = \sum_{(\mathbf{x}_i, \hat{\mathbf{s}}_i, \mathbf{y}_i) \in \mathcal{D}} \alpha(\mathbf{s}, \hat{\mathbf{s}}_i) w(\mathbf{x}, \mathbf{s}, \mathbf{x}_i, \hat{\mathbf{s}}_i) \mathbf{y}_i \quad (9)$$

$$q(\bar{\mathbf{y}}|\mathbf{x}, \mathbf{s}, \mathcal{D}) = 1 - \sum_{(\mathbf{x}_i, \hat{\mathbf{s}}_i, \mathbf{y}_i) \in \mathcal{D}} \alpha(\mathbf{s}, \hat{\mathbf{s}}_i) w(\mathbf{x}, \mathbf{s}, \mathbf{x}_i, \hat{\mathbf{s}}_i) \quad (10)$$

where

$$w(\mathbf{x}, \mathbf{s}, \mathbf{x}_i, \hat{\mathbf{s}}_i) = \frac{\exp(\mathbf{g}(\mathbf{x}, \mathbf{s})^\top \mathbf{g}(\mathbf{x}_i, \hat{\mathbf{s}}_i))}{\sum_{(\mathbf{x}_i, \hat{\mathbf{s}}_i, \mathbf{y}_i) \in \mathcal{D}} \exp(\mathbf{g}(\mathbf{x}, \mathbf{s})^\top \mathbf{g}(\mathbf{x}_i, \hat{\mathbf{s}}_i))} \quad (11)$$

α indicates the function to verify the coincidence of the linguistic explanation and the attribute assigned to each sample. In our setting, we set as $\alpha(\mathbf{s}, \hat{\mathbf{s}}) = \sum_i [[\mathbf{s}[i] = \hat{\mathbf{s}}[i]]]$ where $[[\cdot]]$ is an indicator function of 1 if the condition inside bracket are satisfied, otherwise 0. Note $\alpha(\mathbf{s}, \hat{\mathbf{s}}) \in \{0, 1\}$ as $\|\mathbf{s}\|_0 = 1$. w measures the weight of each referred sample used for prediction. The probability of the sample being assigned to each class is determined to utilize the samples in \mathcal{D} , which match to the linguistic explanation as in (9). An additional ‘‘unknown’’ category $\bar{\mathbf{y}}$ is introduced for convenience, indicating the inability to predict from the input explanations. The remaining weight is assigned to the probability of the ‘‘unknown’’ class, as in (10). In (11), $\mathbf{g}(\mathbf{x}, \mathbf{s})$ is the feature embedding implemented by the neural network as in the right-most in Fig. 3, and the similarity is computed by the dot product in that space following normalization by the softmax function.

While the reasoner attempts to make a decision based on the given explanations, the other two models are trained collaboratively to generate explanations such that the reasoner can reach the appropriate conclusion.

3.6. Training and Inference

We parameterize the joint distribution as in (1), and the lower bound of the objective (2) calculated by (4) and (6) is optimized with regard to the parameters of the neural network models representing $p(\mathbf{s}|\mathbf{x}, \mathbf{y})$, $p(\mathcal{D}|\mathbf{x}, \mathbf{y}, \mathbf{s})$, and $q(\mathbf{y}|\mathbf{x}, \mathbf{s}, \mathcal{D})$. Assuming that the calculation of the expectation over \mathbf{s} is feasible, although that over \mathcal{D} is not, we optimized the model of the selector by sampling, and that of the explainer was optimized directly.

The processing flow in each iteration is as follows:

1. \mathbf{x} is sampled randomly from the training dataset,
2. \mathbf{y} is sampled randomly from the predictor $p(\mathbf{y}|\mathbf{x})$,
3. $p(\mathbf{s}|\mathbf{x}, \mathbf{y})$ is computed for possible \mathbf{s} ,
4. \mathcal{D} is sampled randomly from the selector $p(\mathcal{D}|\mathbf{x}, \mathbf{y}, \mathbf{s})$ for each \mathbf{s} ,
5. For each sampled $(\mathbf{x}, \mathbf{y}, \mathbf{s}, \mathcal{D})$, the objective is calculated by (6) and (4), and the gradient of it w.r.t the weights of all parametric models are computed.
6. All weights are updated by stochastic gradient decent (SGD).

The inference is performed by sequentially sampling variables from the distributions given input \mathbf{x} . When generating linguistic explanations, M identical attribute type is selected whose output value of the predictor is the largest, where M is the number of attribute-examples pairs. used for explanation. For estimating the attribute value, the one

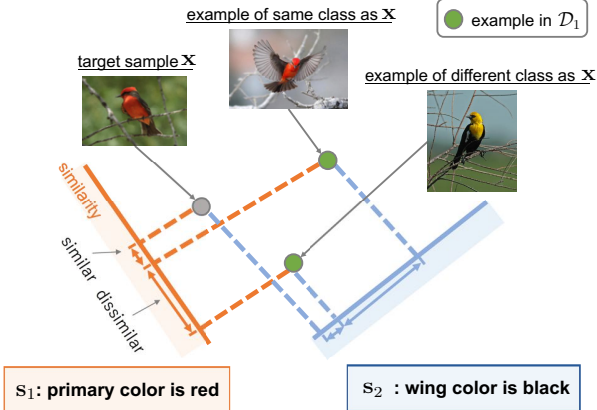


Figure 4: Intuitive understanding of complementary explanations. The reasoner predicts the target sample x (written as gray circles) by referring other samples based on the similarity space (orange and blue) corresponding to each linguistic explanation s_1, s_2 . Considering two pairs of possible explanations (s_1, \mathcal{D}_1) and (s_2, \mathcal{D}_2) , the expected \mathcal{D}_1 (written as green circle) is the one by which the reasoner can reach the correct conclusion with s_1 ; however, this cannot be achieved with s_2 .

that most explains the prediction the best will be selected. In other words, s_1, s_2, \dots having the same attribute type, the value maximizing $q(y|x, s, \mathcal{D})$ is outputted after the corresponding $\mathcal{D}_1, \mathcal{D}_2, \dots$ are sampled from the selector.

3.7. Which explanation is complementary?

By analyzing the proposed method, it provides an intuitive understanding of complementary explanations, from the viewpoint of maximizing interaction information.

To understand which set \mathcal{D} is preferred, we consider (6) where \mathcal{D} relates. Inside the expectation in this equation, the numerator is the output of the reasoner, and the denominator is that averaged over s' . Given x, y , and s , the situation where the ratio becomes large is when the reasoner can reach the correct conclusion y for given linguistic explanation s with \mathcal{D} but it can not when \mathcal{D} is used with other s' . In other words, an example-based explanation is complementary to its linguistic counterpart when it is a *discriminative set* of examples for not only the target but also the linguistic explanation.

The concept of “a set is discriminative” is clearly different from “single example is discriminative” in our framework. This can be understood intuitively by Fig. 4. A reasoner contains a different similarity space for each linguistic explanation s . Here, we consider two possible explanations s_1, s_2 , and \mathcal{D}_1 which is the counterpart of s_1 . In this situation, the desired \mathcal{D} for linguistic explanation s is that the correct class is predicted for the given s , but a wrong one is predicted for different s' . Therefore, the set should contain both of the same/different classes from the predicted

dataset	acc (predictor)	acc (reasoner)	consistency
AADB	0.647	0.646	0.738
CUB	0.694	0.434	0.598

Table 1: The accuracy of identifying the target category of the predictor (target) and reasoner (explain), and the consistency between them.

one. As shown, a naive method of example selection, such as one selecting only one nearest sample from the target, is not appropriate for selecting a complementary explanation.

Considering s , it relates to both terms in (2). For (6), the same claim as that mentioned above can be applied: a complementary linguistic explanation s for examples \mathcal{D} is one where a specific set \mathcal{D} can be derived, instead of another see \mathcal{D}' . As for (4), it can be regarded as a regularizer to avoid weighing excessively on the attribute that can identify the target class without considering examples for selectings.

3.8. Relationship with other methods

The existing works for visual explanation explanations (e.g., [12]) trains the classifier as well as the explanation generator to guarantee that the generated explanations are discriminative for the target class. In this work, we also train the auxiliary classifier (i.e., reasoner) similar to the existing methods; however, it naturally appears in the context of interaction information (mutual information) maximization. Conversely, we found that such an intuitive idea in these works is justified from the information theoretical viewpoint. Similarly, our method shares the idea with methods for generating referring expression (e.g., [28]) in that they utilize auxiliary models.

4. Experiment

We conducted experiments to verify that the proposed method can generate the appropriate explanation. Given a target sample x , our system generates a prediction y from the predictor, and explanations $(s_1, \mathcal{D}_1), (s_2, \mathcal{D}_2), \dots, (s_M, \mathcal{D}_M)$ from explanation models. We evaluated the proposed method by quantifying the properties that the generated explanation should satisfy: (a) fidelity and (b) complementarity. Related to (a), we consider two types of fidelity as follows. (a1) The target value y to be explained should be obtained from the explanations. Moreover, (a2) The outputted linguistic explanation s should be correct. In addition, as for (b), we would like to assess whether the output explanations (s, \mathcal{D}) are complementary to each other. In the following subsections, we describe the evaluation method as well as discussing the obtained results after elucidating the setting of the experiments in subsection 4.1.

4.1. Experimental setting

Dataset In our experiments, we utilized Caltech-UCSD Birds-200-2011 Dataset (CUB) [27] and Aesthetics with

dataset	baseline (random)	baseline (predict)	ours
AADB	0.200	0.572	0.582
CUB	0.125	0.428	0.436

Table 2: The accuracy of identifying the attribute value of our model and that of baselines: selecting attribute value randomly (random), and predicting attributes by the perceptron (predict).

Attributes Database (AADB) [17], both of which hold attributes assigned for all contained images. CUB is a standard dataset for fine-grained image recognition, and it contains 11,788 images in total and 200 categories of bird species. It contains 27 types of attributes, such as “wing pattern” or “throat color.” AADB is a dataset created for the automatic image aesthetics rating. It contains 10,000 images in total and the aesthetic score in [-1.0, 1.0] is assigned for each image. We treat the binary classification by regarding images having non-negative scores as samples of the positive class, and remaining samples as the negative class. Attributes are also assigned as the continuous values in [-1.0, 1.0], and we discretized them according to the range that it belongs to: [-1.0, -0.4), [-0.4, -0.2), [-0.2, 0.2), [0.2, 0.4), or [0.4, 1.0]. It contains eleven types of attributes, including “color harmony” and “symmetry.” Unlike the standard split, we utilized 60% of the test set for CUB, and 30% of the train set for AADB as the candidates of examples \mathcal{X} .

Although CUB dataset is for the fine-grained task, where the inner-class variance of the appearance is considered small, that of AADB is large owing to the subjective nature of the task. We selected these two datasets to assess the influence of the variation of the samples within the class.

Detailed setting To prepare a predictor, we fine-tuned a deep residual network [11] having 18 layers for each dataset, which is pre-trained on ImageNet dataset [6]. The optimization was performed with SGD. The learning rate, weight decay, momentum, and batch size were set to 0.01, 10^{-4} , 0.9, and 64, respectively. When training explanation models, all networks were optimized with SGD without momentum with learning rate 10^{-3} , weight decay 10^{-3} , and batch size 64 for AADB and 20 for CUB. We set k , the number of examples used for explanations, to 10 in all experiments.

Empirically, we found that the linguistic explainer $p(\mathbf{s}|\mathbf{x}, \mathbf{y})$ tended to assign a high probability (almost 1) on only one type of attribute, and small probability (almost 0) to the others. To avoid it, we added an extra entropy term $H(\mathbf{s}|\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{s}} -p(\mathbf{s}|\mathbf{x}, \mathbf{y}) \log p(\mathbf{s}|\mathbf{x}, \mathbf{y})$ to the maximized objective function as our goal is to generate multiple outputs. The implementation was performed on the Pytorch framework [22].

	ours	w/o x	w/o y	w/o s
accuracy	0.646	0.627	0.569	0.613
consistency	0.738	0.689	0.600	0.620

	ours	w/o x	w/o y	w/o s
accuracy	0.434	0.354	0.02	0.153
consistency	0.598	0.492	0.02	0.201

Table 3: The ablation study for the accuracy of identifying the target category on AADB dataset (above) and CUB dataset (below).

4.2. Fidelity

One important factor of the explanation model is the fidelity to the target to be explained. We conducted an experiment to investigate whether the target can be obtained from its explanation. Interestingly, our framework holds two types of paths to the decision. One is the target predictor $p(\mathbf{y}|\mathbf{x})$ to be explained. The other is the route via explanations interpretable for humans, i.e., $\mathbf{y} \rightarrow \mathbf{s}, \mathcal{D} \rightarrow \mathbf{y}'$ through the explainer $p(\mathbf{s}|\mathbf{x}, \mathbf{y})$, selector $p(\mathcal{D}|\mathbf{x}, \mathbf{y}, \mathbf{s})$, and reasoner $q(\mathbf{y}'|\mathbf{x}, \mathbf{s}, \mathcal{D})$. We evaluated the fidelity to the model by the consistency between the interpretable decisions from the latter process and that from the target predictor. In the Table. 1, we reports the consistency as well as the mean accuracy of each models. As shown, the explanation model (written as reasoner) achieved the similar performance as the target model (written as predictor), and considerably high consistency on both datasets.

We also conducted the ablation study to clarify the influence of three variables $\mathbf{x}, \mathbf{y}, \mathbf{s}$ on the quality of explanations. We measured the accuracy in the same manner as above except that we dropped one variable by replacing the vector filled by 0 when generating explanations. The results in Table 3 exhibits that our models put the most importance on the category label out of three. These results are reasonable because it contains information for which the explanation should be discriminative.

The other important aspect for the fidelity in our task is the correctness of the linguistic explanation. In particular, the attribute value (e.g., “red” or “blue” for attribute type “color”) is also estimated during the inference. We evaluated the validity by comparing the predicted attributes with that of grand-truth on the test set. The attribute value that explains the output of the predictor \mathbf{y} the best will be selected as written in subsection 3.6. As baselines, we employed the three layers perceptron with a hidden layer of 512 dimensions (predict). It was separately trained for each type of attribute with SGD. Moreover, we also report the performance when the attribute is randomly selected (random). We measured the accuracy and the results are shown in the Table 2. As shown, our method, which generates linguistic explanations through selecting examples, can predict it as accurate as the direct estimation, and the accuracy is much better than the random selection.

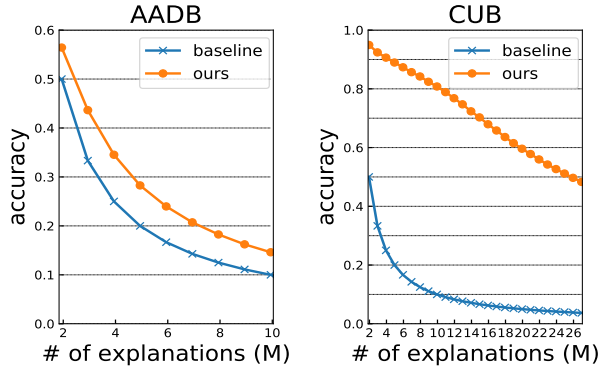


Figure 5: The mean accuracy of identifying the linguistic explanation from the examples on AADB (left) and CUB (right) dataset. The y-axis and x-axis indicates the accuracy and the number of generated explanations.

4.3. Complementarity

To quantify the complementarity of explanations, we investigate how the example-based explanation \mathcal{D} renders the linguistic explanation s identifiable. Specifically, utilizing the reasoner $q(\mathbf{y}|\mathbf{x}, s, \mathcal{D})$, which is trained to reason the target from the explanation, we confirmed whether it can reason to the correct conclusion only from the generated explanation pair as discussed in subsection 3.7. For generated pairs of explanations $(s_1, \mathcal{D}_1), (s_2, \mathcal{D}_2), \dots, (s_M, \mathcal{D}_M)$ of which attribute type is identical, we computed the output of the reasoner as $q_{ij} = q(\mathbf{y}|s_i, \mathcal{D}_j)$ ($1 \leq i, j \leq M$) for \mathbf{y} obtained from the predictor. Selecting the index having the maximum value for all j as $i^* = \arg \max_i q_{ij}$, we verified $i^* = j$. The mean accuracy is compared with a baseline that outputs the same examples for all s_i and results are shown in Fig. 5. The x-axis of the figure indicates the number of the generated explanations (i.e., M).

On both datasets, the accuracy of our model is better than the baseline. Furthermore, as shown in Fig. 6, we observe that the diagonal element has a high value on the confusion matrix. These results demonstrate the ability of our method to generate complementary explanations. The difference in the performance between the proposed method and baseline on AADB is lower than on CUB. We conjecture that one reason is the difference between appearance and attributes. AADB dataset contains highly-semantic attributes (e.g., “color harmony”) compared with those in CUB (e.g., “color” or “shape”). Such the semantic gap may hinder to construct the discriminative set which renders the attribute identifiable.

4.4. Output example

An output of our system on CUB dataset when the number of explanations is two is shown in Fig. 7. In this example, the combination of linguistic and example-based explanation seems compatible, where it will not make sense

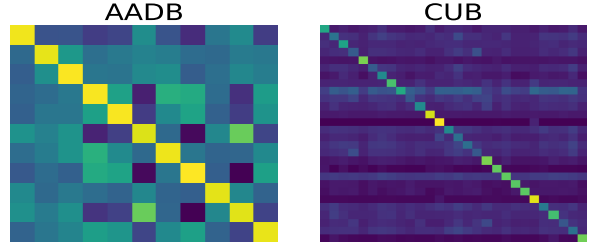


Figure 6: The confusion matrix of identifying the attribute type from the examples on AADB (left) and CUB (right) dataset.

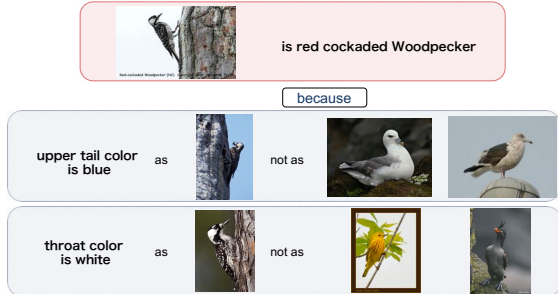


Figure 7: Example output of our system on CUB dataset.

if these pairs are switched. For instance, the below linguistic explanation “throat color is white” may not be a good explanation for the above examples.

Although not the primary scope in this work, the proposed task may be extended to machine teaching task, where the machine *teaches* to human by showing examples iteratively.

5. Conclusion

In this work, we performed a novel task, that is, generating visual explanations with linguistic and visual examples that are complementary to each other. We proposed to parameterize the joint probability of variables to explain, and to be explained by the three neural networks. To explicitly treat the complementarity, auxiliary models responsible for the explanations were trained simultaneously to maximize the approximated lower bound of the interaction information. We empirically demonstrated the effectiveness of the method by the experiments conducted on the two visual recognition datasets.

6. Acknowledgement

This work was partially supported by JST CREST Grant Number JPMJCR1403, Japan, and partially supported by the Ministry of Education, Culture, Sports, Science and Technology (MEXT) as “Seminal Issue on Post-K Computer.” Authors would like to thank Hiroharu Kato, Toshihiko Matsuura for helpful discussions.

References

- [1] Oxford living dictionaries. Oxford University Press. complement <https://en.oxforddictionaries.com/definition/complement>. 1
- [2] L. Anne Hendricks, R. Hu, T. Darrell, and Z. Akata. Grounding visual explanations. In *ECCV*, 2018. 2
- [3] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015. 2
- [4] J. Chen, L. Song, M. J. Wainwright, and M. I. Jordan. Learning to explain: An information-theoretic perspective on model interpretation. In *ICML*, 2018. 2, 4
- [5] P. Dabkowski and Y. Gal. Real time image saliency for black box classifiers. In *NIPS*, 2017. 2
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 7
- [7] C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. Efros. What makes paris look like paris? *ACM Transactions on Graphics*, 31(4), 2012. 2
- [8] R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *ICCV*, 2017. 2
- [9] S. Gao, G. Ver Steeg, and A. Galstyan. Variational information maximization for feature selection. In *NIPS*, 2016. 4
- [10] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 2
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 7
- [12] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell. Generating visual explanations. In *ECCV*, 2016. 2, 6
- [13] A. Jain, A. Gupta, M. Rodriguez, and L. S. Davis. Representing videos using mid-level discriminative patches. In *CVPR*, 2013. 2
- [14] E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. 4
- [15] A. Kanehira, K. Takemoto, S. Inayoshi, and T. Harada. Multimodal explanations by predicting counterfactuality in videos. In *CVPR*, 2019. 2
- [16] A. Kanehira, L. Van Gool, Y. Ushiku, and T. Harada. Viewpoint-aware video summarization. In *CVPR*, 2018. 2
- [17] S. Kong, X. Shen, Z. Lin, R. Mech, and C. Fowlkes. Photo aesthetics ranking network with attributes and content adaptation. In *ECCV*, 2016. 7
- [18] O. Mac Aodha, S. Su, Y. Chen, P. Perona, and Y. Yue. Teaching categories to human learners with visual explanations. In *CVPR*, 2018. 2
- [19] C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016. 4
- [20] W. McGill. Multivariate information transmission. *Transactions of the IRE Professional Group on Information Theory*, 4(4):93–111, 1954. 2, 3
- [21] D. H. Park, L. A. Hendricks, Z. Akata, A. Rohrbach, B. Schiele, T. Darrell, and M. Rohrbach. Multimodal explanations: Justifying decisions and pointing to the evidence. In *CVPR*, 2018. 2
- [22] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017. 7
- [23] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017. 2
- [24] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013. 2
- [25] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery of mid-level discriminative patches. In *ECCV*. 2012. 2
- [26] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *NIPS*, 2016. 5
- [27] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical report, 2011. 6
- [28] L. Yu, H. Tan, M. Bansal, and T. L. Berg. A joint speaker-listener-reinforcer model for referring expressions. In *CVPR*, 2017. 6
- [29] J. Zhang, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff. Top-down neural attention by excitation backprop. In *ECCV*, 2016. 2
- [30] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. 2
- [31] B. Zhou, Y. Sun, D. Bau, and A. Torralba. Interpretable basis decomposition for visual explanation. In *ECCV*, 2018. 2