

Bachelorarbeit

# **Uncertainty in Recurrent Decision Tree Classifiers**

Eberhard Karls Universität Tübingen  
Mathematisch-Naturwissenschaftliche Fakultät  
Wilhelm-Schickard-Institut für Informatik  
Explainable Machine Learning  
Stefan Wezel, [stefan.wezel@student.uni-tuebingen.de](mailto:stefan.wezel@student.uni-tuebingen.de), 2020

Bearbeitungszeitraum: von-bis

Betreuer/Gutachter: Prof. Dr. Zeynep Akata, Universität Tübingen  
Betreuer: Stephan Alaniz, Max Planck Institut für Informatik

# Abstract

Recurrent Decision Tree Classifiers have proven to be capable of providing explanations of their classification while yielding state of the art results in prediction accuracy on several image classification tasks. However, we show that they may utilize features that they are highly uncertain about. We hypothesize that when using uncertainty information, the RDTc can provide more faithful explanations and become more applicable in real-life scenarios. Based on RDTc, we propose two models that utilize uncertainty information. We either enhance the vocabulary of RDTc with an uncertainty token, allowing for a ternary decision tree, or, we restrict the RDTc from using uncertain attributes. We investigate how uncertainty information can be used in generating interpretable model outputs and how it affects the models performance on several benchmark tasks.

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Related Work</b>	<b>7</b>
2.1	Neural Networks . . . . .	7
2.1.1	Historical Context . . . . .	7
2.1.2	Multilayer Perceptron . . . . .	8
2.1.3	Recurrent Neural Networks . . . . .	8
2.1.4	Convolutional Neural Networks . . . . .	9
2.2	Explainable Machine Learning . . . . .	9
2.2.1	Decision Trees . . . . .	9
2.2.2	Leveraging Attributes . . . . .	10
2.2.3	Multi-agent Communication . . . . .	10
2.3	Uncertainty . . . . .	11
2.3.1	Bayesian Inference . . . . .	12
2.3.2	Variational Inference . . . . .	12
2.3.3	Gaussian Processes . . . . .	13
2.3.4	Modeling Uncertainty in Neural Networks . . . . .	13
<b>3</b>	<b>Uncertain RDT<sub>C</sub></b>	<b>15</b>
3.1	RDT <sub>C</sub> . . . . .	15
3.1.1	Communication . . . . .	15
3.1.2	AbL . . . . .	16
3.1.3	RDT . . . . .	16
3.1.4	Learning . . . . .	17
3.2	Uncertainty . . . . .	18
3.2.1	Estimating Epistemic Uncertainty with Dropout Neural Networks	18
3.2.2	Getting Uncertainty Information . . . . .	21
3.2.3	Removing Uncertain Attributes . . . . .	21
3.2.4	Extending the Vocabulary . . . . .	21
3.2.5	Random Attribute Removal . . . . .	22
<b>4</b>	<b>Experiments</b>	<b>24</b>
4.1	Datasets . . . . .	24
4.2	Statistics on Uncertainty . . . . .	25
4.3	Testing on Out-of-domain Data . . . . .	25
4.4	Results . . . . .	26
4.4.1	Comparing RDT <sub>C</sub> to Other Models . . . . .	27

## Contents

4.4.2 Comparing Different RDTC Models . . . . .	27
<b>5 Discussion</b>	<b>29</b>
5.1 Reduced Performance with Uncertainty . . . . .	29
5.2 The Right Kind of Uncertainty? . . . . .	29
5.2.1 Estimating Heteroscedastic Uncertainty . . . . .	30
5.2.2 Risk versus Uncertainty . . . . .	30
5.2.3 Beyond Attribute Uncertainty . . . . .	31
5.3 A Practical Example . . . . .	31
5.4 Carbon Footprint Estimation . . . . .	31
<b>6 Conclusion</b>	<b>33</b>

# Chapter 1

## Introduction

The recent surge in popularity of machine learning methods has made the fields ambivalent nature obvious. Public opinion is split between hype and mistrust. Systems, using machine learning have failed at seemingly simple tasks [32], displayed racist and sexist behavior [23, 68, 77] and even caused fatalities [76]. Such cases, where an artificial intelligence has made severe mistakes are funny at best and outright dangerous at worst. Most importantly, they have made it clear that, despite being deployed more frequently than ever before, machine learning has yet to prove itself [68, 14, 77]. These cases show that methods that allow introspection and behave faithfully are more important than ever.

Addressing such issues, Alaniz and Akata [2] propose a model that can reliably justify its classification through stating utilized features in a decision tree. Here, we build on this method. We extend the model, so that it is aware of and able to express its own uncertainties. We hypothesize that such a model may be perceived by a user as more trustful than a model that does not.

We, as humans have developed an intuition how to make decisions under uncertainty [5] and are able to operate in a domain of limited and noisy data. Moreover, uncertainty is an important measure in science that allows to express how well something is known or can be known. Scientists from fields like medicine, physics or meteorology rely on uncertainty information for making and communicating decisions or predictions. Outside the scientific world, uncertainty information is important in finance or law among others [43, 9].

Uncertainty information is valuable. Especially in real-world scenarios, where we might encounter only limited, or noisy data to learn from. Where predictions, far into the future are required. Or, where decisions with high consequences have to be made. As machine learning becomes more prevalent in our everyday lives, in science, finance, or law, the need for methods, utilizing uncertainty becomes more evident.

Unfortunately, the output of most neural network architectures does not contain uncertainty information. For example, in classification, typical architectures use a softmax layer, to turn model outputs into a probability density function [53] over classes. Since it is a probability distribution, the softmax output is often misinterpreted as a measure of how certain a model is regarding each class [69]. However, this distribution only arises from each value relative to every other value and is a point estimate, turned into a distribution. The point estimate cannot carry any uncertainty information and the softmax certainly does not add such. This allows a softmax class probability to be overconfident for a highly uncertain class [18]. Our key contributions are to propose a model, that allows to (1) study the effect

## Chapter 1. Introduction

of (epistemic) uncertainty for attribute based recurrent decision tree classifier and actively uses and expresses uncertainty information. We (2) test the proposed models performance on popular datasets and compare it to different state-of-the-art methods. We investigate uncertainties in the model and data and evaluate its affect on performance.

Such a model may be applied in settings, where a high confidence classification for out-of-domain examples may have severe consequences. For example in medical settings, in law or in autonomous driving. There, in cases of high uncertainty, a human could be notified and consulted, preventing the model from making a (potentially wrong) high-risk decision. The model would still provide a decision tree as explanation, which could be utilized by a human user to make a final decision.

# Chapter 2

## Related Work

Building on, often centuries old, mathematical foundations, and on more recent psychological and neurobiological insights, machine learning has evolved throughout the last decades and a vast corpus of literature has emerged from this development. Here we will briefly go over some historic aspects, introduce important concepts and give an overview of the current state of explainable methods and uncertainty in machine learning.

### 2.1 Neural Networks

Artificial neural networks (ANNs) have become an integral part in the lives of many. Inspired by the cognitive processes of the mammalian brain, several architectures and learning algorithms have been proposed.

#### 2.1.1 Historical Context

The idea of artificial intelligence (AI) is dating back to antiquity and was debated in various domains ranging from chemistry [54] to philosophy and mathematics [42]. A watershed moment in its development is marked by McCulloch and Pitts [46], who modeled a neuron in an approach that unified findings from neurophysiology and mathematics. Later, Hebb [26] proposed a theory of how the strengthening and weakening of connections between neurons encodes gained knowledge(?). This idea was also applied in the perceptron [64] which could learn distinguish linearly separable classes. Marvin and Seymour [45] showed that the perceptron was more useful when connecting many and using a middle layer in a so called multilayer perceptron (MLP). If this hidden layer has a non-linear activation function it could learn to classify not linearly separable data. This idea of many simple but connected units remains a key principle in ANNs. However, simulating the behavior of neural networks requires many computations and the sequential nature of hardware available at the time prohibited connectionist ideas to be applied in practice.

Over the decades, ANNs have risen and fallen in popularity among AI researchers. Different architectures and learning algorithms have been proposed but many failed to establish. Important milestones were Rumelhart et al. [65], who extended learning rules to MLPs in the backpropagation algorithm, the Long short-term memory by Hochreiter and

## Chapter 2. Related Work

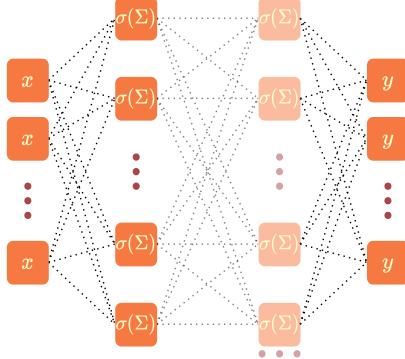


Figure 2.1: A MLP can consists of an input layer, an arbitrary number of hidden layers, and an output layer.

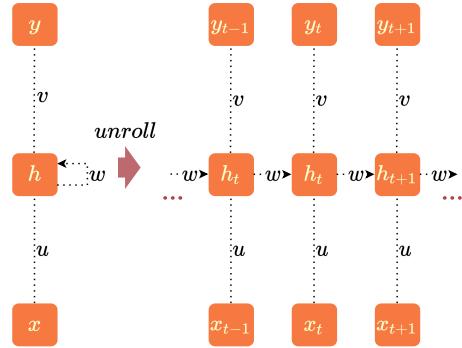


Figure 2.2: The recurrent neural network takes in an input for every time step as well as its own hidden state.

Schmidhuber [27], and the convolutional neural net (CNN), proposed by LeCun et al. [41] which allowed to classify handwritten digits.

However, these methods still failed to establish far beyond the realm of research. The development of graphics processing units (GPUs) and the availability of large scale datasets have caused this to change in the 2010s where for some tasks ANNs have exceeded human performance [66].

### 2.1.2 Multilayer Perceptron

The MLP is a rather simple architecture. It consists of an input layer, an arbitrary number of hidden layers and an output layer. Each neuron is typically connected to each neuron in the next layer. Thus each neuron gets the weighted output from the last layer's neurons as input. Non-linear activation functions allow the network to learn not linearly separable problems.

### 2.1.3 Recurrent Neural Networks

Many types of data don't provide samples of fixed length. This is the case for most language or time series data among others. A MLP, however, is incapable of processing such samples due to its fixed amount of input neurons. To alleviate this, recurrent neural networks (RNNs) use the output of their neurons as an input after being multiplied with a weight matrix. Many concepts of RNNs exist [20]. Here we will have a closer look at a very simple example, depicted in Figure 2.2. The RNN takes in an input  $in_t$  which is transformed by weight matrix  $u$  and added to the hidden state of the last time step times weight matrix  $w$ . This forms the hidden state of the current time step. To obtain an output at a given time step, the hidden state is multiplied with weight matrix  $v$ . The update rules can be stated as:

$$h_t = f_{act}(u \circ in_t + w \circ h_{t-1}) \quad (2.1)$$

$$out_t = f_{act}(v \circ h_t) \quad (2.2)$$

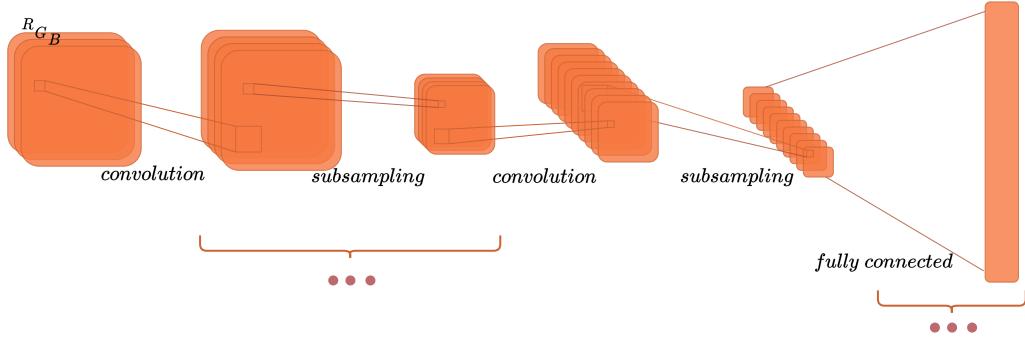


Figure 2.3: A CNN with two convolutional and subsampling layers which can be repeated arbitrary times.

A RNN can be trained by the Backpropagation through time algorithm (BPTT). Werbos [73] propose to unfold a RNN for each time step or element of the input. This unfolding makes it equivalent to a feed-forward ANN (such as the MLP introduced in Subsection 2.1.2) and weights can be adjusted according to the gradients.

## 2.1.4 Convolutional Neural Networks

CNNs have proven to be useful tool for various image related tasks. They are a type of feed-forward network with the underlying concept of learning filters that correspond to features in the data. Typically this is done over multiple layers with the intention of learning to recognize low level concepts such as edges in layers close to the input layer and high level concepts in layers further away from the input. In practice, after convolutions and non-linearities are applied, the results are usually subsampled, using a pooling operation to reduce the amount of necessary parameters.

## 2.2 Explainable Machine Learning

While interpretable models exist in machine learning, they are often outperformed by modern ANNs with a large number of parameters, on a variety of tasks [21, 22, 67, 59]. Understanding the processes inside a model with a large number of parameters is not trivial. Thus, various approaches have been developed to understand a models decision process for given data. Such methods, that are interpretable or can explain their decisions are inevitable for safety critical applications. They allow a human user to understand how a model came to a certain output and allow for a faithful human-AI interaction to emerge.

### 2.2.1 Decision Trees

Decisions trees provide an inherently interpretable model. Given data with features  $a_1, \dots, a_n$  a decision tree iteratively splits data according to conditions for each attribute. Splits occur for each attribute  $a_i$  where samples from data are either above or below a given threshold for their values of  $a_i$ . Optimal splits can be found by learning algorithms or by expert

## Chapter 2. Related Work

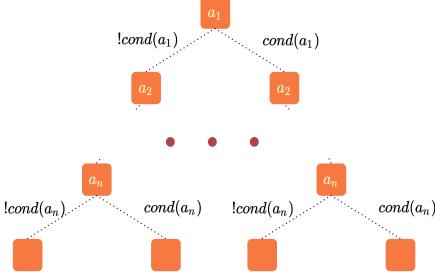


Figure 2.4: A decision tree splits data according to values of attributes.

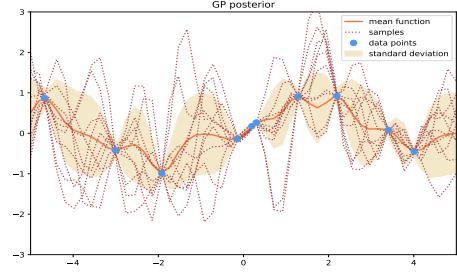


Figure 2.5: A Gaussian process is a distribution over the space of functions that describe data (blue).

knowledge. Nodes represent the resulting splits and edges represent the corresponding answer. Decision trees are a frequently used tool in data mining [74]. Popular algorithms to generate decision trees from data include Iterative Dichotomiser 3 and C4.5, both suggested by Quinlan [60, 61]. Another popular algorithm is Classification and Regression Trees, introduced by Breiman et al. [7]. The key concept behind these algorithms is to create splits that maximize information gain for each available attribute.

In order to reduce size and increase interpretability of generated trees, pruning can be applied [48]. Pruning can also reduce the risk of overfitting ([31]).

### 2.2.2 Leveraging Attributes

While decision trees are suitable for tabular data, applying them on image data is not trivial. Images are parameterized by pixels rather than attributes, so attributes have to be provided as side information. Ideally, those attributes should correspond to features that characterize classes found in the image data. Therefore, a popular way of acquiring side information is by humans label the data according to characteristics suggested by domain experts. Wah et al. [72] let users label different kind of bird species according to attributes and collect them in the CUB dataset. Another dataset containing attributes as side information is proposed by Xian et al. [75].

Attributes can be used to create an interpretable output of a neural network. Akata et al. [1] propose to learn a linear mapping between feature representations learned by a neural network and learned attribute representations. In other areas, attributes are used as well. Kulkarni et al. [37] and Ordonez et al. [55] use attributes for image captioning. Lampert et al. [39] and Palatucci et al. [57] use attributes for zero-shot learning.

### 2.2.3 Multi-agent Communication

Communication between different agents is a popular technique in reinforcement learning [24, 40, 8, 29, 11]. Foerster et al. [16] propose to let one agent send messages containing categorical symbols which can be used by a second agent to solve a problem. Multi-agent communication has also been used for the purpose of creating interpretable outputs by Rodriguez et al. [63] who take in account that different agents may have differing

### 2.3. Uncertainty

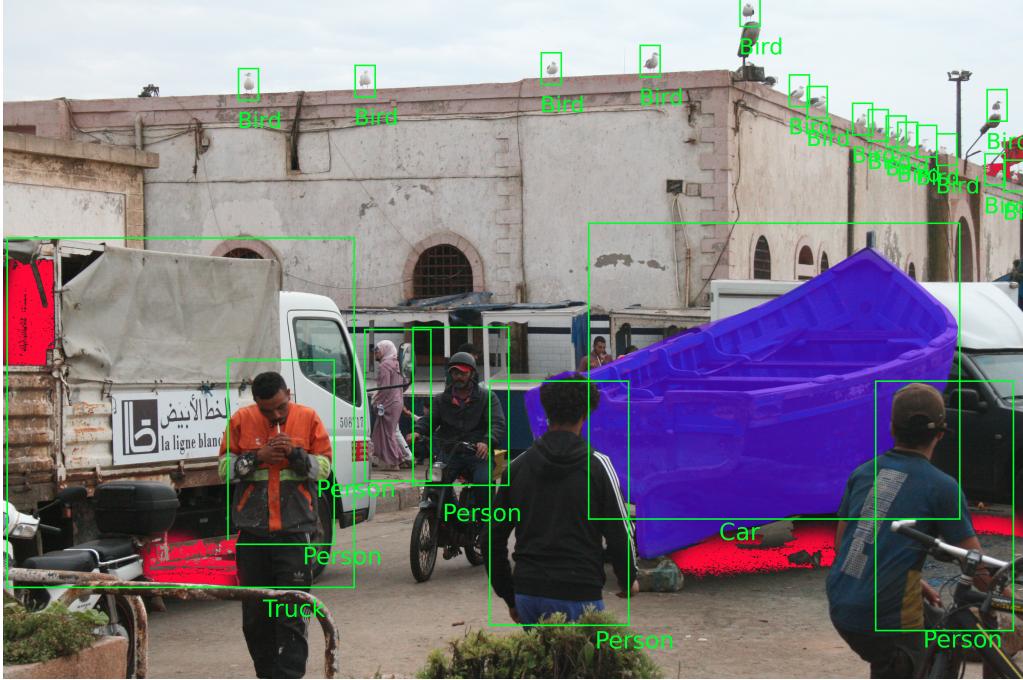


Figure 2.6: Areas with high aleatoric uncertainty are marked red. These are mostly due to occlusions in the image. The area in purple is due to high epistemic uncertainty, since the model has never seen a boat before. It misclassifies it as a car but is uncertain about this classification

understanding of concepts. Communication between two agents is also the basis for learning in the RDTC, proposed by Alaniz and Akata [2], which our work builds on.

## 2.3 Uncertainty

Uncertainty is an important factor in human decision making. Modeling uncertainty in machine learning methods is useful for ensuring responsible decision making of algorithms and allows a user to build trust in these models. Uncertainty can arise from data. This is called aleatoric uncertainty. It compromises all kinds of noise inherent to the data. If this noise is static (i.e. limited precision of a sensor) it is called homoscedastic uncertainty. If noise is not static (i.e. occlusions in an image) the uncertainty is called heteroscedastic uncertainty.

On the other hand, a model can also be uncertain about its own parameters. This uncertainty is called epistemic uncertainty and is high if a model has seen few data. By giving it more data, it can be reduced. Those two kinds of uncertainty are displayed for a computer vision task in Figure 2.6.

## Chapter 2. Related Work

### 2.3.1 Bayesian Inference

Applying Bayes' theorem to update prior beliefs as more data becomes available is called Bayesian inference [6]. Bayes' theorem was introduced by Bayes [3] and is stated as

$$p(a|b) = \frac{p(b|a) \cdot p(a)}{p(b)}.$$

Intuitively put, it computes the posterior conditional probability of a random variable  $a$  given, the likelihood of  $b$  being true, given  $a$  is true:  $p(b|a)$ , a prior belief  $p(a)$  and a marginalization term  $p(b)$ . Methods of Bayesian inference are popular in different fields of science, such as neurology, psychology, medicine, epidemiology, geography and machine learning [17, 71, 34, 58, 13, 19]. They can be used to incorporate prior knowledge, multiple sources of evidence, or, in general, construct large joint probability density functions [70].

### 2.3.2 Variational Inference

Methods of variational inference are used when intractable integrals are encountered in Bayesian inference [30]. Suppose, we want to find the predictive posterior  $p(y^*|x^*, X, Y)$  for some datapoint  $x^*$ , we need to compute

$$p(y^*|x^*, X, Y) = \int p(y^*|x^*, w)p(w|X, Y)dw.$$

The true posterior over model parameters, given data  $p(w|X, Y)$  is not trivial to find and is often intractable in practice. Therefore an approximate distribution  $q_\theta(w)$ , with optimizable parameters  $\theta$  is used. Ideally,  $q_\theta(w)$  should be similar to the true posterior. Thus, we minimize the KL-divergence between the approximate distribution and the true posterior.

$$KL(q_\theta(w)||p(w|X, Y)).$$

Given an approximate distribution, the predictive posterior can be written as

$$p(y^*|x^*) \approx \int p(y^*|x^*, w)q_\theta(w)dw.$$

In order to find the variational distribution  $q_\theta^*(w)$  with the optimal set of parameters, we have the objective

$$\mathcal{L}_{VI} := \int q_\theta(w) \log p(Y|X, w) dw - KL(q_\theta(w)||p(w))$$

that is optimized in respect to parameters  $\theta$ . Here,  $p(w)$  is a prior that is often assumed to be a standard Gaussian. According to Bishop [4], maximizing the log evidence lower bound is equivalent to minimizing the KL-divergence but does not require having access to the true posterior  $p(w)$ . This effectively leaves us with an optimization task where we need to compute derivatives instead of integrals. This, in turn, allows us to apply powerful automatic differentiation tools, such as deep learning frameworks.

## 2.3. Uncertainty

### 2.3.3 Gaussian Processes

Gaussian processes are a powerful, generative model that can be applied to a variety of machine learning tasks [44]. They build on the notion that data can be described by (possibly infinitely) many functions. Gaussian processes (GPs) assign probabilities to functions that describe given data. Starting with a prior distribution over functions a GP returns a posterior over the space of functions, given new data. Of course, distributions over functions are somewhat hard to define. Therefore, we will view a function  $f(x)$  as a vector where the  $x$  denotes the index that will retrieve the function value at this point. Now, a GP is a joint normal distribution over the points of this vector, defined by a mean function  $\mu(x)$  and a covariance function  $\Sigma(x)$  with  $\Sigma(x) = \Sigma_{ij} = k(x_i, x_j)$ . The kernel function  $k$  has to return a positive definite matrix and should return a high output for similar  $x_i$  and  $x_j$ .

GPs thus view each  $x_i$  as a random variable, that is dependent on all other  $x_j \neq x_i$ .

To perform inference in a GP, one simply takes the conditional probability density distribution (PDF) of a given  $x_i$  conditioned on all other  $x_j \neq x_i$ . Since Gaussian distributions are closed under conditioning, the resulting PDF is Gaussian as well.

To retrieve a meaningful posterior from a GP, it has to be trained in a supervised fashion. We will describe the necessary steps in pseudocode.

```

-create a prior-
 $x_{test} = \{x_n\}$ 
 $K_{prior} = \text{kernel}(x_{test}, x_{test}, w)$ 
 $L_{prior} = \text{cholesky\_decomposition}(K_{prior})$ 
 $f_{prior} = L_{prior} \cdot \mathcal{N}(\text{dim} = (n, num\_samples); 0, 1)$ 
-take new data into account and compute new posterior-
 $x_{train}, y_{train} = \text{training data}$ 
 $K = \text{kernel}(x, x, w)$ 
 $L = \text{cholesky\_decomposition}(K)$ 
 $K_s = \text{kernel}(x_{train}, x_{test}, w)$ 
 $L_k = \text{solve}(L, K_s)$ 
 $\mu = L_k^T \cdot \text{solve}(L, y_{train})$ 
-draw samples from posterior-
 $L = \text{cholesky\_decomposition}(K_{ss} - L_k^T \cdot L_k)$ 
 $f_{posterior} = \mu + L \cdot \mathcal{N}(\text{dim} = (n, num\_samples); 0, 1)$ 
- prediction at test point  $x_i$  is now the sample's value at position  $i$ 
```

### 2.3.4 Modeling Uncertainty in Neural Networks

Deep neural networks (DNNs) are applied in safety critical processes where information about a model's uncertainty is of crucial importance [50], [36]. Standard DNN's however do not yield uncertainty information. Bayesian neural networks (BNNs) are an architecture that can model uncertainty reliably [49]. However, they require modeling a probability density function over each parameter which causes a high computational cost. A method using variance arising from dropout has been proposed by Gal and Ghahramani [18] who proofed dropout in neural networks is equal to Monte-Carlo integration in GPs and thus

## Chapter 2. Related Work

can be used to model epistemic uncertainty.

With a GP, we want to find a posterior distribution over function that explain given data. We can view a finite model like a neural network as an approximation to a GP [10]. Thus, by optimizing a model, we minimize the Kullbach-Leibler (KL) divergence between our finite model and the corresponding GP [18]. In practice, however, the GP's posterior often requires computing intractable integrals. To alleviate this, methods of variational inference, such as Monte-Carlo integration are applied. Monte-Carlo integration in a GP corresponds to averaging forward passes in a dropout neural net. According to Gal and Ghahramani [18], the resulting variance can be views as the model's uncertainty. Additionally, Kendall and Gal [33] propose to combine aleatoric and epistemic uncertainty in a unified approach. For classification tasks, this has been simplified by Kwon et al. [38].

## Chapter 3

# Uncertain RDT

The Recurrent decision tree model, proposed by Alaniz and Akata [2], consists of two communication agents. Their common goal is to solve a classification task through communication. The Recurrent Decision Tree (RDT) agent has no direct access to data and can only ask question regarding the absence, presence or uncertainty of attributes. Given answers by the Attribute-based Learner (AbL) agent it can then store the answers thus incrementally building a decision tree that is used for the final classification. The learned decision three allows introspection.

Based on their work, we propose a model, where the AbL is able to express its uncertainty for given attributes. We use this uncertainty information by either preventing the RDT from using uncertain attributes or by allowing the AbL to answer with 'I don't know' additionally to the options 'Yes' or 'No'<sup>1</sup>.

### 3.1 RDT

Before we introduce our modifications to the RDT model, we will give a detailed explanation of its architecture, learning algorithm and how the two agents communicate.

#### 3.1.1 Communication

The two agents communicate to each other through discrete values. For a given image, the AbL returns a an answer  $\hat{a}$  that is a tensor with binary values with size number of attributes  $\times$  decision size. The intention behind using only binary values in the AbL's answer is increased interpretability. Each row in this tensor indicates the presence/absence of an attribute and the columns correspond to the decision options 'Yes' or 'No'. Strictly speaking, the AbL thus answers all possible questions the RDT can ask, in advance to any questions. The RDT is then able to ask questions through posing an index  $c_t$ . By accessing the AbL's answer at this index it receives answer  $d_t = \hat{a}[d_t]$ . This information then used to build an internal decision tree, used for a classification output and to come

---

<sup>1</sup>We provide a Pytorch implementation of our model at [github.com/wastedsummer/urdtc](https://github.com/wastedsummer/urdtc)

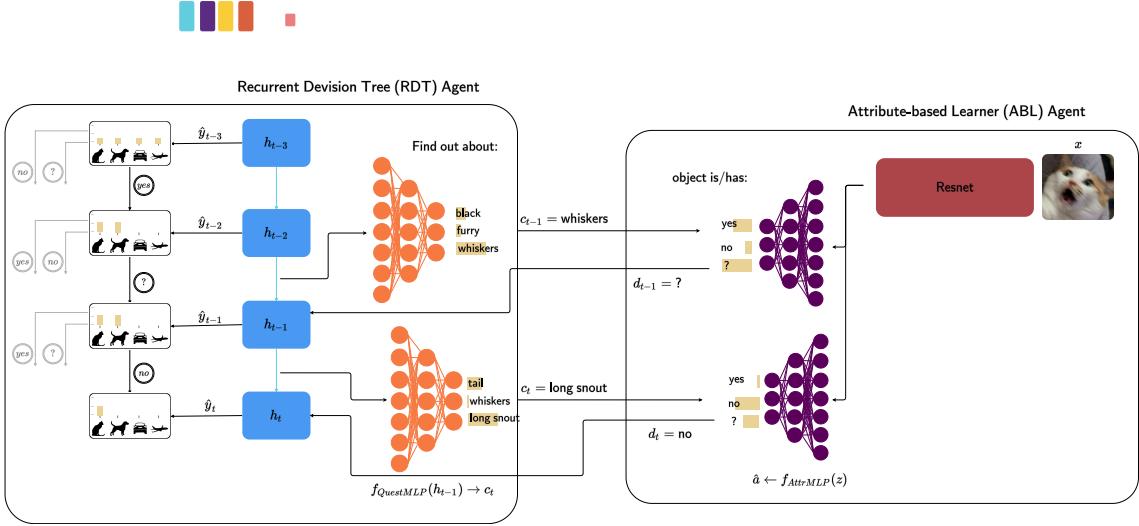


Figure 3.1: The RDT agent asks questions about presence, absence or uncertainty of attributes. The answers, given by the AbL agent are used by the RDT agent to make a classification each iteration.

up with a new question. In the following, we will go over each agent in more detail and explain how this model is trained.

### 3.1.2 AbL

The AbL agent is a vision model, consisting of a CNN and a MLP. The CNN learns image features and the MLP learns a mapping from features to human-annotated attributes (if available). For all experiments, we use a ResNet [25] to extract features. For a given image  $x$ , the CNN extracts feature vector  $z$ . This is passed to the  $f_{AttributeMLP}$  and the resulting vector of logits is put through a softmax. The output is a tensor  $p(\hat{a}|z)$  with shape  $number\_attributes \times decision\_size$ . This corresponds to a probability distribution over the decision option for each attribute. In order to get an interpretable, binary response, in the forward pass, the result of the softmax is put through an arg max function. This is, however, not differentiable. In the backward pass it is therefore replaced by the identity function. Thus, the final response of the AbL agent for a given image  $x$  is  $\hat{a} = argmax(p(\hat{a}|z))$ . The RDT agent can then access answers for specific attributes by indexing.

### 3.1.3 RDT

The RDT agent consists of several models. At its core lies an Long-short Term Memory (LSTM). This recurrent model allows the RDT to come up with new questions, based on it's hidden state  $h_{t-1}$  from the last (or initial) communication step.

The  $f_{QuestMLP}$  is responsible for thinking of new questions by returning an index that can be used to access attribute information in the response from  $f_{AttributeMLP}$ . The output of  $f_{QuestMLP}$  is a categorical distribution  $\log p(c_t|h_{t-1})$  over possible indices  $c_t \in \{1, 2, \dots, |A|\}$ . In order to get an actual index, the final output needs to be discrete. Thus we need to sample

### 3.1. RDTC

from this categorical distribution. This is done via the Gumbel softmax function:

$$GumbelSoftmax(\log \pi) = \frac{\exp((\log \pi_i + g_i)/\tau)}{\sum_{j=1}^K \exp((\log \pi_j + g_j)/\tau)}, \text{ where} \quad (3.1)$$

$$\log p_i = \log p(c_t|h_{t-1}), \text{ and thus:} \quad (3.2)$$

$$c_t = GumbelSoftmax(\log p(c_t|h_{t-1})). \quad (3.3)$$

where for  $\tau \rightarrow 0$ , the distribution approaches a one-hot vector with a 1 at the maximum of the categorical distribution and 0 everywhere else.

The resulting scalar  $c_t$  then can be used to index  $\hat{a}$ . The response is  $d_t$  which indicates absence, presence or uncertainty of the attribute at index  $c_t$ . Question  $c_t$  and answer  $d_t$  then update explicit memory  $\mathcal{M}^{(t)} = \mathcal{M}^{(t-1)} \oplus (c_t, d_t)$ . Note that during test time, rather than sampling  $c_t$ , the index with the highest probability is chosen.

The explicit memory  $\mathcal{M}^{(t)}$ ,  $c_t$ ,  $d_t$  and the LSTM's last hidden state  $h_{t-1}$  are then used to computer the current hidden state  $h_t$ , which is then, in turn, used by the  $f_{QuestMLP}$  again to pose a new question.

The explicit memory  $\mathcal{M}^{(t)}$  is also the basis for the class prediction. An additional MLP, the  $f_{ClassMLP}$  returns class probabilities  $\hat{y}_t$  based on  $\mathcal{M}^{(t)}$  for every communication step:

$$\hat{y}_t = f_{ClassMLP}(\mathcal{M}^{(t)}). \quad (3.4)$$

The classification is the the final output of the model and is used to calculate a loss that can be optimized.

#### 3.1.4 Learning

Ideally, the model should come to a classification in as few communication steps as possible. A tree resulting from few communication steps is shallow and easier to interpret. To keep the model from using deep trees, the classification loss is divided by the number of time steps (communication steps):

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \mathcal{L}_{CE}(y, \hat{y}_t) = -\frac{1}{T} \sum_{t=1}^T \sum_i y_i \cdot \log \hat{y}_{t,i}. \quad (3.5)$$

This ensures, that the lower the number of time steps  $T$ , the lower the whole term and thus the lower the loss.

Besides the classification loss, an additional attribute loss can be leveraged in order to encourage learning attributes that correspond to provided human annotated side information.

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T [(1-\lambda) \mathcal{L}_{CE}(y, \hat{y}_t) + \lambda \mathcal{L}_{CE}(\alpha_{y,c_t}, \hat{\alpha}_{c_t})] \quad (3.6)$$

How much each loss contributes to the total loss can be determined by hyperparameter  $\lambda$ .

## 3.2 Uncertainty

Uncertainty information is a crucial part of robust vision systems. In safety critical application, it could be avoided that the model uses high-risk strategies by leveraging attributes it is highly uncertain about for its classification. Or, if all attributes are deemed highly uncertain (for example on completely ood data), the model could decide to consult a human instead of making an overconfident classification.

Uncertainty comes either from variance in data (aleatoric uncertainty) or from variance in model parameters (epistemic uncertainty). Thus, there is an inherent relationship between variance and uncertainty. The variance of a random variable, which is in our case the variational predictive distribution  $p(\hat{y}|x)$  with  $y \sim p(\hat{y}|x)$ , is defined as:

$$Var_q(p(\hat{y}|x)) = \mathbb{E}_q[(y - \mathbb{E}[y])^2] \quad (3.7)$$

$$= \mathbb{E}_q[yy^T] - \mathbb{E}_q[y]\mathbb{E}_q[y]^T. \quad (3.8)$$

To decompose the term into variance arising from data and variance from the model  $q$  itself, according to Kwon et al. [38], we can write it as:

$$Var_q(p(\hat{y}|x)) = \mathbb{E}_q[yy^T] - \mathbb{E}_q[y]\mathbb{E}_q[y]^T \quad (3.9)$$

$$= \underbrace{\int_{\Omega} [diag(\mathbb{E}_{p(\hat{y}|x,w)}[\hat{y}]) - \mathbb{E}_{p(\hat{y}|x,w)}[\hat{y}] \cdot \mathbb{E}_{p(\hat{y}|x,w)}[\hat{y}]] q_{\theta}(w) dw}_{\text{aleatoric}} \quad (3.10)$$

$$+ \underbrace{\int_{\Omega} [\mathbb{E}_{p(\hat{y}|x,w)}[\hat{y}] - \mathbb{E}_{q_{\theta}(\hat{y}|x,w)}[\hat{y}]] [\mathbb{E}_{p(\hat{y}|x,w)}[\hat{y}] - \mathbb{E}_{q_{\theta}(\hat{y}|x,w)}[\hat{y}]]^T}_{\text{epistemic}}. \quad (3.11)$$

### 3.2.1 Estimating Epistemic Uncertainty with Dropout Neural Networks

Epistemic uncertainty corresponds to the variance of the models parameters. In a traditional dropout neural network, there is no random dropping out of neurons in test stage and weights are fix. This would result in no variance of model outputs for multiple forward passes given an input. However, dropping out neurons in a random manner results in an inherent variance of model parameters. Gal and Ghahramani [18] proof that this variance corresponds to epistemic (model) uncertainty. Here, we will go through their proof.

The key idea is that a dropout neural net is equivalent to an approximation of a Gaussian process. To randomly dropout neurons, we sample binary values that are either 0 or 1 for every neuron. The sampled value is then multiplied with the neurons value. In the following, the dropout neural network with  $L$  layers the weight matrices are denoted by  $W_{l \in L}$ , the bias vectors by  $b_l$ , the network output given datapoint  $x_i \in X$  with size  $N$  by  $\hat{y}_i$ , and the true value by  $y_i \in Y$ . The objective of the network is to minimize a given loss function  $E(y_i, \hat{y}_i)$ . Additionally, a weighted regularization term can be added to the loss function. A popular technique is to use  $L_2$  regularization, weighted by a weight decay term  $\lambda$ . The final resulting loss term is:

$$\mathcal{L}_{\text{dropout}} := \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{y}_i) + \lambda \sum_{l=1}^L (\|W_l\|_2^2 + \|b_l\|_2^2).$$

### 3.2. Uncertainty

For the special case of a single hidden layer dropout network, we can rephrase this objective as

$$\mathcal{L}_{dropout} := \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{y}_i) + \lambda (\|W_1\|_2^2 + \|W_2\|_2^2 + \|b\|_2^2) \quad (3.12)$$

, where  $W_1$  is the weight matrix connecting the input to the hidden layer, and  $W_2$  connecting the hidden to the output layer. The bias  $b$ . Let us take a step back to a very general point of view. A problem, we typically encounter in machine learning is how to predict a function value  $y^*$  from a datapoint  $x^*$ . In Bayesian terms, we want to find the predictive posterior

$$p(y^*|x^*, X, Y) = \int p(y^*|x^*, w)p(w|X, Y)dw \quad (3.13)$$

where  $w$  is a set of adjustable parameters. In order to model our function of interest well, we need to find optimal values for  $w$ . If we break down Equation 3.13 further, we get

$$p(y^*|x^*, X, Y) = \int \underbrace{p(y^*|x^*, w)}_{\mathcal{N}(\underbrace{y; \hat{y}(x, w)}_{\hat{y}(x, w = \{W_1, W_2, \dots, W_n\})}, \tau^{-1} \mathbf{I}_D)} \cdot \underbrace{p(w|X, Y)dw}_{intractable} \quad (3.14)$$

, where we see that the posterior over  $w$ , given data is intractable. Therefore we need to approximate it with variational distribution  $q_\theta(w)$ . We want this distribution to mimic the true posterior closely. Therefore, we want to find

$$\operatorname{argmin}_\theta KL(q_\theta(w) \| p(w|X, Y))$$

As we have seen in Section 2.3.2, this would require access to the true posterior, which we cannot compute. Therefore we maximize the log evidence lower bound

$$\mathcal{L}_{VI} := \int q_\theta(w) \log p(Y|X, w) dw - KL(q_\theta(w) \| p(w))$$

and assume  $p(w)$  to be a standard Gaussian. As discussed also in Section 2.3.2, we effectively turn computing integrals into an optimization problem, where we need to compute derivatives. Using the optimized variational distribution  $q_\theta^*(w)$ , we can compute an approximate predictive posterior

$$q_\theta(y^*|x^*) = \int p(y^*|x^*, w) q_\theta^*(w) dw$$

for a given datapoint  $x^*$ . This in turn can be approximated with

$$q_\theta(y^*|x^*) = \sum_{t=1}^T p(y^*|x^*, w_t)$$

where each  $w_t$  is a sample from our variational distribution  $p_\theta^*(w)$ .

Now, in Gaussian processes, as discussed in Section 2.3.3, we need a covariance function of the form

$$K(x, y) = \int \mathcal{N}(w; 0, l^{-2} I_Q) p(b) \sigma(w^T x + b) \text{sigma}(w^T y + b) dw db$$

### Chapter 3. Uncertain RDT

to get a covariance matrix. To avoid computing integrals, this covariance function can be approximated with

$$\hat{K}(x, y) = \frac{1}{K} \sum_{k=1}^{k=1} \sigma(w_k^T x + b_k) \sigma(w_k^T y + b_k)$$

where  $w_k$  is sampled from  $\mathcal{N}(0, l^{-2} I_Q)$  and  $b_k$  from  $p(b)$ . This is called Monte-Carlo integration [47].

This gives us following predictive posterior for our GP:

$$\begin{aligned} w_k &\sim \mathcal{N}(0, l^{-2} I_Q), w_d \sim \mathcal{N}(0, l^{-2} I_K), b_k \sim p(b) \\ W_1 &= [w_k]_{k=1}^K, W_2 = [w_d]_{d=1}^D, b = [b_k]_{k=1}^K, w = \{W_1, W_2, b\} \\ p(y^*|x^*, w) &= \mathcal{N}\left(y^*; \sqrt{\frac{1}{K}\sigma(x^*W_1 + b)W_2}, \tau^{-1}I_N\right) \end{aligned}$$

,and finally:

$$p(y^*|x^*, X, Y) = \int p(y^*|x^*, w)p(w|X, Y)dw.$$

We use  $q_\theta(w) = q_\theta(W_1)q_\theta(W_2)q_\theta(b)$  as an approximation to  $p(w|X, Y)$ , where

$$q_\theta(W_1) = \Pi q_\theta(w_q), \text{ and } q_\theta(w_q) = p_1 \mathcal{N}(m_q, s^2, I_K) + (1-p_1) \mathcal{N}(0, s^2, I_K), \quad (3.15)$$

$$\text{with } p_1 \text{ as a probability } \in [0, 1], s > 0, \text{ and } M_1 = [m_q]_{q=1}^Q \in \mathbb{R}^{K \times D} \quad (3.16)$$

We eventually want to recover a dropout neural network (as specified in Eq. 3.12) from our GP. In our dropout network we rely on sampling from a Bernoulli random variable. In our GP, we, so far, sample from a mixture of Gaussian distributions. We cannot change this, because the KL-divergence between a discrete distribution (Bernoulli) and a continuous distribution (Gaussian) is not defined. Therefore we will make two assumptions about the Gaussian mixture in Eq. 3.15.

We set  $p$  to a fix value. Also, we want the second term of the sum to be zero. To sample zero from a zero-centered Gaussian distribution with a very high probability, we simply set the standard deviation to a very small value. Small enough, so that a device with limited precision will view it as zero and always return zero. These two assumptions leave us with a Gaussian distribution that resembles a Bernoulli distribution. Thus, we are effectively sampling from a Bernoulli random variable in our GP, with a still defined KL-divergence. We get

$$\mathcal{N}_{GP-MC} \approx \log p(Y|X, \hat{w}) - \frac{p_1 l^2}{2} \|M_1\|_2^2 - \frac{p_2 l^2}{2} \|M_2\|_2^2 - \frac{l^2}{2} \|m\|_2^2, \text{ with } \hat{w} \sim q_\theta(w)$$

as approximated log evidence lower bound. To recover the dropout objective from Eq. 3.12, we can scale it by a constant  $\frac{1}{N\tau}$ .

$$\mathcal{L}_{GP-MC} \propto -\frac{1}{2N} \sum_{n=1}^N \|y_n - \hat{y}_n\|_2^2 - \frac{p_1 l^2}{2N\tau} \|M_1\|_2^2 - \frac{p_2 l^2}{2N\tau} \|M_2\|_2^2 - \frac{l^2}{2N\tau} \|m\|_2^2$$

and set hyperparameters  $\tau$  and  $l$  to appropriate values.

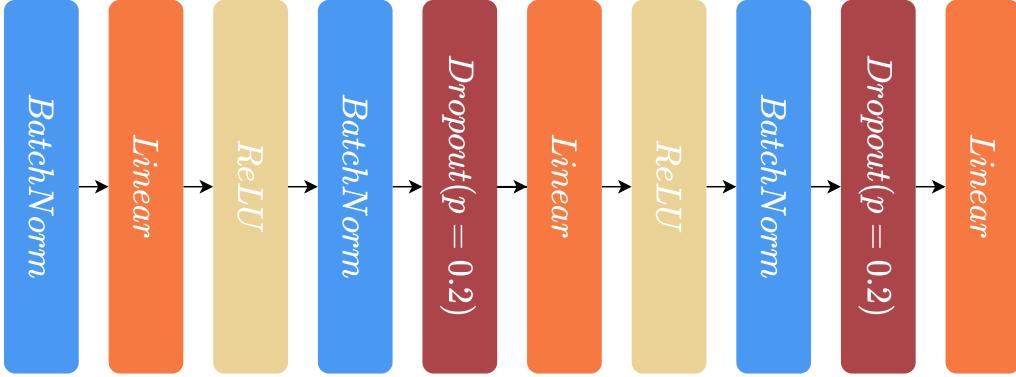


Figure 3.2: We used both BatchNorm and dropout layers for  $f_{attrMLP}$ .

### 3.2.2 Getting Uncertainty Information

We tested different configurations for  $f_{AttributeMLP}$  and using a combination of dropout and batchnorm layers between, fully connected layers and non-linearities worked best. This setup is displayed in Figure 3.2. To retrieve uncertainty information, given the image features  $z$  from the ResNet, we do  $n$  forward passes and compute the standard deviation of the results. As proofed by Gal and Ghahramani [18], this is equal to epistemic uncertainty. So theoretically, for examples, that are far from the distribution, the model was trained on, it should yield high epistemic uncertainty. After retrieving uncertainty information, we replace selection logits from  $f_{QuestMLP}$  that are on indices with high uncertainty by  $-inf$  so they can not be picked by the Gumbel softmax. We call this strategy remove uncertain attributes (remRDTc). In another strategy, uncertain attributes are denoted with a 1 in a extended decision tensor. We refer to as extended vocabulary (extRDTc).

### 3.2.3 Removing Uncertain Attributes

After calculating uncertainty for each attribute, in this strategy, we simply manipulate the logits tensor from  $f_{QuestMLP}$  so that at each index, where a uncertain attribute is, the selection logit tensor has a  $-inf$  value. This ensures that this attribute cannot be picked by the Gumbel softmax. In turn, no uncertain attribute can be picked as index and therefore cannot be posed as a question to the AbL.

In order to prevent the model from misusing uncertainty information, we have to ensure that the model does not get any gradients where uncertain attributes are.

### 3.2.4 Extending the Vocabulary

Here, we resort to an extended vocabulary of answers for  $f_{AttrMLP}$ . Instead of just allowing Yes and No, we also allow the AbL to answer with ?. Thus,  $f_{AttrMLP}(z) \in \{0, 1\}^{num\_attributes \times 3}$ . However, since,  $f_{AttrMLP}(z) \in \{0, 1\}^{num\_attributes \times 2}$  we need to append the uncertainty information to  $d$  while avoiding conflicting answers. This process is shown, for a two attribute example in Figure 3.4. After receiving the output from  $f_{AttrMLP} = d_{init}$  we initiate the uncertain column ? with a vector of 0s and append this to the  $d_{init}$  to create  $d_{temp}$ . After

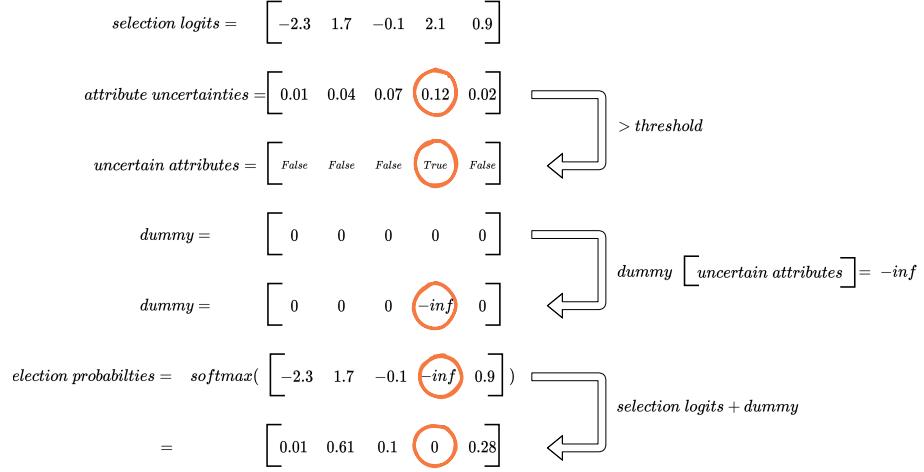


Figure 3.3: In selection logits, the output from  $f_{QuestMLP}$ , we replace values at indices where uncertain attributes are with a  $-inf$  value. This assures that this attribute can never be picked as  $c_t$  by the Gumbel softmax.

computing uncertainty vector  $?$ , we also create a negated version  $!$  of it. This indicates all attributes where the model is certain. Subsequently, those are now used to clear rows. Vector  $!$  has 0s where the model is uncertain. We use this to replace values in rows of uncertain attributes with 0s and get  $d_{cleared}$ . We also initiate a dummy tensor, that is of the same dimension as the  $d_{cleared}$  and is all 0s but in the column to the very right, where it is all 1s. Rows of this tensor are now cleared, using  $?$ , wherever  $?$  has zeros, thus only rows with uncertain attributes remain and have a 1 in the column responsible for uncertainty. Finally, this tensor, and  $d_{cleared}$  are added up to create our final decision  $d$  which can now be indexed by the RDT.

When using this strategy in training, we need to make sure, that no gradients actually come from uncertainty, so it is not misused and maximized in order for the model to profit from the additional column. We need to ensure that gradients only come from attributes, where the model is certain. For the operation, as described above, this is the case. We detach  $?$  and  $!$  from the computational graph. The dummy tensor is a leaf node, so no gradient is computed. Also, where values in  $d$  are 0, their activation's derivative is 0, so gradients are 0. This is the case in rows of uncertain attributes. Which gradients are used and which are blocked in the previously used example is displayed in Figure 3.5.

### 3.2.5 Random Attribute Removal

Similar to the remRDTc strategy, described in Subsection 3.2.3, we prevent the RDT from asking questions. However, instead of blocking indices where uncertain attributes are, we randomly generate a tensor and block indices based on it. The probability of an index being blocked is similar to the probability of an attribute being deemed uncertain to allow fair comparison.

### 3.2. Uncertainty

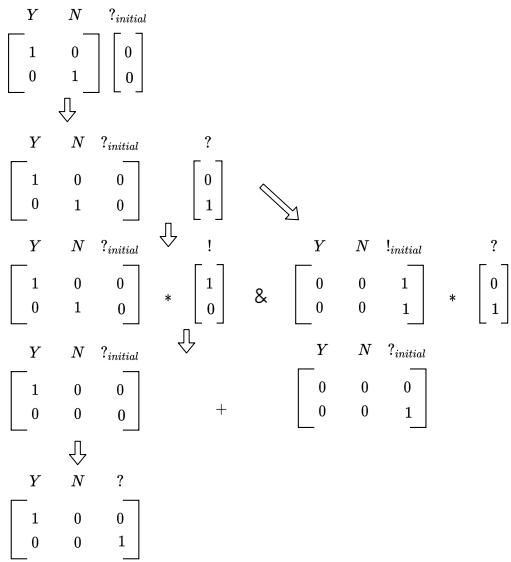


Figure 3.4: The original answer by the ABL is merged with uncertainty information so that the answer for each attribute is either 'yes', 'no' or '?'. Here we show a small example where attribute size is only 2.

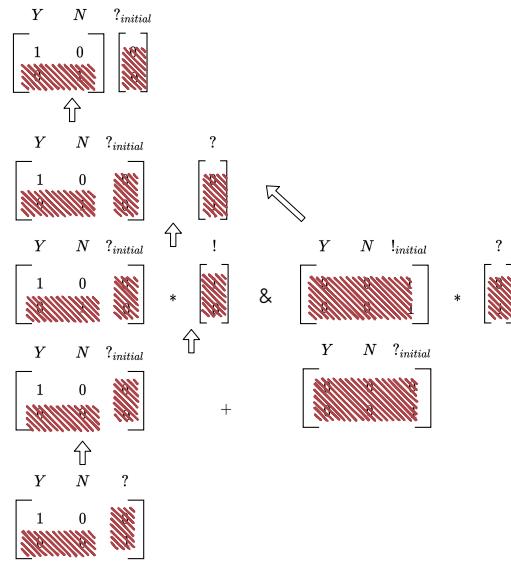


Figure 3.5: Gradients coming from attributes that are deemed uncertain are detached from the computational graph. Other elements in the tensor are either 0 (so are their respective gradients), or stem from leaf nodes.

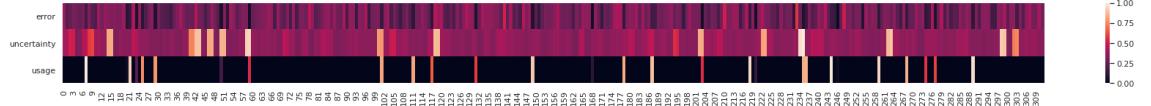


Figure 4.1: Normalized attributes misclassification frequency, uncertainty and usage. Only a small subsets of attributes is actually used by the model.

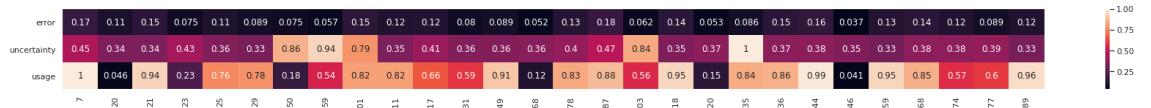


Figure 4.2: Here, we only show normalized attributes misclassification frequency, uncertainty and usage for attributes that are used by the model.

## Chapter 4

# Experiments

We evaluated our proposed model on different image classification datasets which provide attributes as side information. We also investigate the models behavior by quantitatively analyzing uncertainties and their affects.

### 4.1 Datasets

We test our model on image datasets that provide attributes as side information. Popular datasets for this setting are Animals with Attributes 2 (AwA2) [75], aPY [15] and CUB [72]. Datasets with attributes are frequently used for zero- or few-shot learning tasks. In the explainable machine learning community they are popular since attributes provide a human interpretable, natural language parametrization of features.

AwA2 is a set of 37,322 images that can be subdivided into 50 different classes of animals. Human annotators were tasked to label each image, choosing from 85 available attributes that describe each animal class. The aPY dataset provides 15,339 images, annotated with 64 different attributes that describe each of the 32 classes. Finally, the most challenging dataset is CUB. Here, the task is to classify 11,788 different images into 200 different classes. For this task, each images is annotated with 312 descriptive attributes.

## 4.2. Statistics on Uncertainty

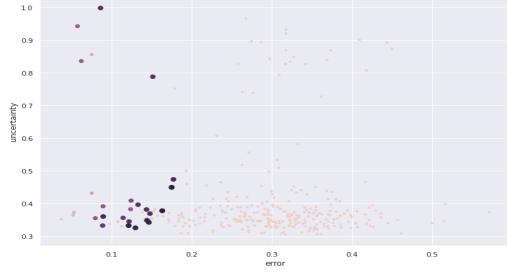


Figure 4.3: ...

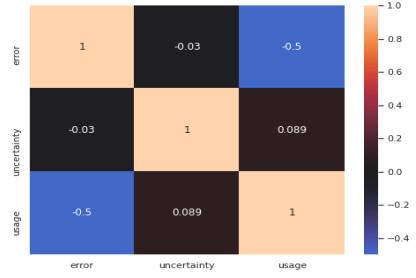


Figure 4.4: ...

## 4.2 Statistics on Uncertainty

The purpose of the proposed model is not only to evaluate it on benchmark datasets but also to quantitatively investigate uncertainties in the model and the data. We examined correlations between variables such as attribute uncertainty, usage and misclassification on the CUB dataset. From Figure 4.1 it becomes obvious that only a small subset off attributes as actually used to make decisions. We can already see, that most of the attributes that are used, however, are rarely misclassified. This is also visible in Figure 4.2 where only attributes that are used in the decision tree are shown. As apparent in the top row, all attributes here are classified correctly quite reliably. However, there seems to be no obvious correlation between uncertainty and misclassification as among the frequently used attributes (and thus reliably classified correctly) there are some attributes that have high uncertainty values.

In Figure 4.3 we can see that this is true and uncertain attributes appear among all levels of attribute accuracy. Figure 4.4 proves further proofs this as the correlation between the two is close to 0. It also shows high negative correlation between attribute misclassification frequency and usage. We can view misclassification frequency as inverse accuracy and thus have a positive correlation of 0.5 between usage and accuracy, suggesting that the model actively avoids attributes that it is likely to classify.

## 4.3 Testing on Out-of-domain Data

An important aspect of our model is that it is not only self aware regarding uncertainty, but also allows for human introspection. Here, we test whether examples, far from the training data, the model has seen cause high uncertainty. We use a attribute-zero-shot setting, where we exclude certain classes from training data, so that some attributes are never seen by the model in training. For our setting, we use the CUB dataset, as it is the most challenging out of the introduced three.

A classical zero-shot setting or a generalized zero shot setting (both described by [75]) would not yield meaningful insights here. As in our proposed model, we only consider epistemic attribute uncertainty, uncertainty of classes should only depend on whether attributes attributes have been seen by the model and not whether classes have seen by the model, a possible weakness of our model, we further discuss in 5.2.3. We want to construct a setting, where certain attributes do not occur in training data, and thus have to find classes, that have exclusive attributes. In CUB, we find a set of 24 classes that posses 5

## Chapter 4. Experiments

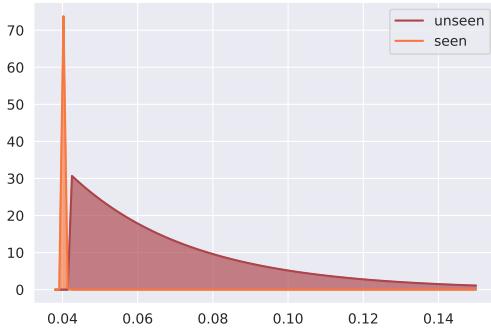


Figure 4.5: Exponential maximum likelihood model for uncertainty per example. Examples from classes that have unseen attributes have higher attribute uncertainties.

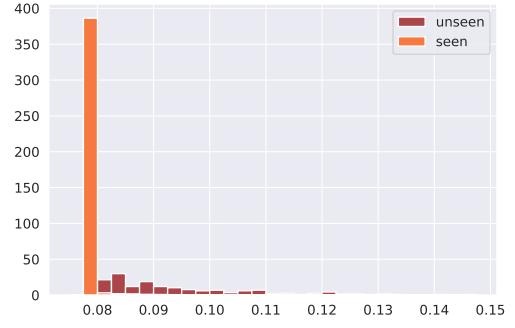


Figure 4.6: Distribution of uncertainty values of seen and unseen attributes in zero-shot-attribute setting (See Section 4.3). The model is more uncertain about attributes it has not encountered in the training set.

attributes that do not occur in any of the other classes<sup>1</sup>.

After training our model on the proposed split, we investigate (1) uncertainties of classes that contain unseen attributes (Figure 4.5) and (2) uncertainties of the unseen attributes themselves (Figure 4.6). For (1), we computed the mean uncertainty per example. For (2), we consider all uncertainty values of each attribute. We observe that the 24 classes, possessing unseen attributes have a overall higher uncertainty. Figure 4.5 displays the modeled distribution of mean uncertainty per example, divided into classes that only posses attributes, occurring in the training data, and classes that posses 5 attributes, never encountered during training. We observe, that the latter have a overall higher uncertainty. Displayed in Figure 4.6, is the distribution of uncertainty values of attributes available during training, and the uncertainty values of attributes, unseen by the model, before the test stage. Here, we also observe an overall higher uncertainty for the latter.

## 4.4 Results

We compare our model to different other approaches, that either exemplify the current state of the art or are explainable.

ResNet-152 [25] is not explainable. Here, the ResNet was pre-trained on ImageNet data [12] data and then fine-tuned on the specific tasks. We also use the ResNet trained in this manner in our proposed uRDT model to extract features. When we train the uRDT model, we keep the weights in the ResNet fixed and only adjust all other parameters in the model. XDT and DT are both traditional decision tree models that serve as a baseline. For DT Alaniz and Akata [2] propose following method: DT is given features extracted by the previously introduced pre-trained ResNet. Splits are created according to values in extracted feature vector. For each split, one feature is chosen and then a split gets created. This is repeated until either each leaf node corresponds to exactly one class or a early stopping-criterion is met. Moreover, Alaniz and Akata [2] introduce XDT, where each

---

<sup>1</sup>This setting can be replicated by using the `zero_attr_train_test_split.txt` file in the provided Github repository.

## 4.4. Results

split is given a semantic meaning, in order for the splits to correspond to interpretable attributes and not only features. The XDT is therefore not given features, directly stemming from the ResNet, but attributes, learned by a MLP that was trained to predict attributes given ResNet features. They train both DT and XDT using the CART algorithm [7]. As splitting criterion, they chose the Gini impurity index over entropy-based methods due to its computational advantage [62].

Another explainable model is Deep Neural Decision Forest (dNDF), proposed by Kortscheder et al. [35]. They learn an ensemble of trees by optimizing routing probabilities in each trees internal nodes. Every path through a tree eventually leads to class distribution, represented in leaf nodes. All internal nodes are sigmoid-activated, differentiable stochastic functions, parameterized by a optimizable  $\theta$ . The final prediction is an average over the resulting decision tree forest. On downside of the resulting averaged stochastic routing decisions is the reduced interpretability as they allow for multiple possible routes. The RDTC models only yield one possible route for a given example. Therefore, we consider only dNDF using a single tree instead of an ensemble mean.

### 4.4.1 Comparing RDTC to Other Models

	AWA2	aPY	CUB
ResNet [25]	$98.2 \pm 0.0$	$85.1 \pm 0.6$	$79.0 \pm 0.2$
DT	$78.0 \pm 0.4$	$64.3 \pm 0.6$	$19.3 \pm 0.3$
dNDF[35]	$97.6 \pm 0.2$	$85.0 \pm 0.6$	$73.8 \pm 0.3$
RDTC[2]	$98.0 \pm 0.1$	$85.7 \pm 0.7$	$78.1 \pm 0.2$
XDT	$73.9 \pm 0.9$	$59.9 \pm 1.5$	$4.9 \pm 1.3$
aRDTC[2]	$98.1 \pm 0.0$	$85.3 \pm 0.3$	$77.9 \pm 0.6$
remRDTC(ours)	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$
extRDTC(ours)	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$

### 4.4.2 Comparing Different RDTC Models

We evaluated different configurations of our model on the three datasets. We compared the original RDTC model to models using our two proposed strategies of attribute removal and extended vocabulary. We also included a baseline, that randomly removes attributes, based on the average probability of an attribute having an uncertainty higher than the specified threshold.

## Chapter 4. Experiments

### **AWA2**

	aRDTC [2]	Random Baseline	remRDTC	extRDTC
Class				
Attribute				

### **aPY**

	aRDTC [2]	Random Baseline	remRDTC	extRDTC
Class				
Attribute				

### **CUB**

	aRDTC [2]	Random Baseline	remRDTC	extRDTC
Class				
Attribute				

# Chapter 5

## Discussion

Uncertainty estimation in deep learning is not trivial. It is an active topic of research and therefore subject of ongoing change and discussion. In Section 4.4, we have seen that models utilizing uncertainty information are outperformed by models which do not. We will critically reflect on our work and discuss where the performance issues might stem from, whether epistemic uncertainty (that is only considering attributes) is useful in our setting and how we could alleviate remaining issues. We also take a look at a hypothetical example, on how uncertainty information could be used.

### 5.1 Reduced Performance with Uncertainty

Both of our proposed models, do not perform as well as the original RDT, proposed by Alaniz and Akata [2]. We hypothesize why this might be the case for (1) our removing attributes model (rRDT) and (2) our extended vocabulary model (uRDT).

In our rRDT model, we actively remove attributes and thus information that otherwise could be utilized by the RDT agent. In Figure 4.2, we see that the model uses attributes it is highly uncertain about. If such attributes cannot be chosen by the RDT, it has less information (albeit risky information) available about an image to make a classification. In our extended vocabulary strategy this is not the case. Here, the RDT always gets information about any attribute it may ask. However, since we avoid using gradients, coming from uncertain attributes, we still might loose information. This can lead to slower convergence time. It also might lead to reduced performance since the path chosen, by the gradient descent algorithm might differ from that chosen by a model, with additional gradients. Thus, it may find a different local minimum.

### 5.2 The Right Kind of Uncertainty?

In our approach, we only considered epistemic uncertainty, which is the uncertainty arising from variance in the model's parameters. However, in closed datasets, as the ones considered, epistemic uncertainty will not change very much after the model has seen every training example. This happens within one epoch.

### 5.2.1 Estimating Heteroscedastic Uncertainty

Heteroscedastic is a function of data and thus can be learned. Kendall and Gal [33] propose to learn an additional output  $\sigma^2$  and a loss function:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{\|y_i - f(x_i)\|^2}{2\sigma(x_i)^2} + \frac{1}{2} \log \sigma(x_i)^2. \quad (5.1)$$

If a prediction is wrong, it encourages a high value for  $\sigma^2$ , thus lowering its impact. In a variational inference setting (which our dropout model is corresponding to) for classification with sample size  $T$ , Kendall and Gal [33] propose to model aleatoric uncertainty as:

$$\frac{1}{T} \sum_{t=1}^T \text{diag}(\hat{\delta}_t^2) \quad (5.2)$$

Alternatively, a different approach was proposed by Kwon et al. [38] that does not require learning  $\sigma^2$ .

$$\frac{1}{T} \sum_{t=1}^T \text{diag}(\hat{p}_t) - \hat{p}_t \hat{p}_t^T \quad (5.3)$$

where  $\hat{p}_t = \text{softmax}(\pi)$  and  $\pi$  are the unscaled logits. This is an approximation to the variational posterior:

$$\int_{\Omega} \left[ \text{diag}(\mathbb{E}_{p(y^*|x^*, w)}[y^*]) - \mathbb{E}_{p(y^*|x^*, w)}[y^*] \cdot \mathbb{E}_{p(y^*|x^*, w)}[y^*] \right] q_{\theta}(w) dw, \quad (5.4)$$

which is the diagonal of the expected outputs, subtracted by the product of the expected outputs with themselves and finally multiplied with the variational posterior over the weights  $w \in \Omega$ . For  $T \rightarrow \infty$ , equation 5.3 will converge to equation 5.4.

### 5.2.2 Risk versus Uncertainty

Osband [56] suggests to distinct the variance arising from dropout neural networks, from actual uncertainty. He criticizes that the approximated posterior, which the mean and variance from multiple forward passes in a dropout network are equal to, according to Gal and Ghahramani [18], do not concentrate with more data available. According to him, this variance is rather associated with what is called risk in decision theory. There, the term risk refers to the inherent stochasticity of a model. However, this is what we previously introduced as epistemic uncertainty. This issue is thus rather about semanticity. We therefore acknowledge it but stick to the convention of differentiating between epistemic and aleatoric uncertainty.

### 5.3. A Practical Example

#### 5.2.3 Beyond Attribute Uncertainty

As the only point, where uncertainty information is retrieved in our model is in the  $f_{AttrMLP}$  which maps ResNet features to attributes, our model can only consider attribute uncertainty. Thus, as soon as an attribute is encountered in training, uncertainty decreases. This may be unhelpful in realistic scenarios, where an unknown class is encountered that can be described by a set of known attributes as such an example would not yield high uncertainty. Typical zero-shot settings would be such scenarios [75]. Replacing  $f_{ClassMLP}$  with a Dropout MLP and computing uncertainty there as well may potentially alleviate this issue and could allow further insights.

## 5.3 A Practical Example

We demonstrated a model that can tell a user reliably about its uncertainties in an explainable decision process. In a realistic setting this could provide the user with valuable insights. Instead of only yielding the final classification and an explanation tree, the model can also tell the user how uncertain it is about each part of the tree. If the user sees that the model was uncertain about several attributes it used, classification could be done manually.

Let's consider an example (displayed in Figure 5.1) where an ornithologist is tasked to do a comprehensive survey of bird species and their respective numbers. Since our ornithologist is not given any helper, she decides to use a drone and a computer vision software that shoots a picture of every bird it encounters (1). After letting the drone fly around in the survey area for some days, the ornithologist inspects the data and finds that the drone collected so much data that it hardly can be examined manually (2). Luckily, she has heard of our UncertainRDT and decides to use it on the data. For each image, the model yields a classification and an explanation. However, due to climate change some bird species that typically don't occur in the survey area have made their way into collected data. The model has never seen such birds that typically appear only much more south and thus is highly uncertain and lets the ornithologist know (3). The ornithologist sees that the model is unaware of this species and therefore cannot reliably classify it. However, since she knows of this species she can quickly classify it manually.

## 5.4 Carbon Footprint Estimation

To produce our results, we relied on GPUs provided by the University of Tübingen ML Cluster of Excellence and Google. We are aware of the profound impact, of the CO<sub>2</sub> emissions caused by producing the required electricity, on the world's climate. The GPUs provided by Google were Nvidia Tesla K80 which operate up to the thermal design power of 300 W [51]. Devices of the University of Tübingen GPU cluster are Tesla V100-SXM2 which, at full capacity, operate up to a thermal design power of 300 W [52]. These specifications are without considering cooling.

While no single model was trained for more than 50 epochs, if we consider all models and all iterations of models, we accumulate about 2000 epochs of training. Training time for

## Chapter 5. Discussion

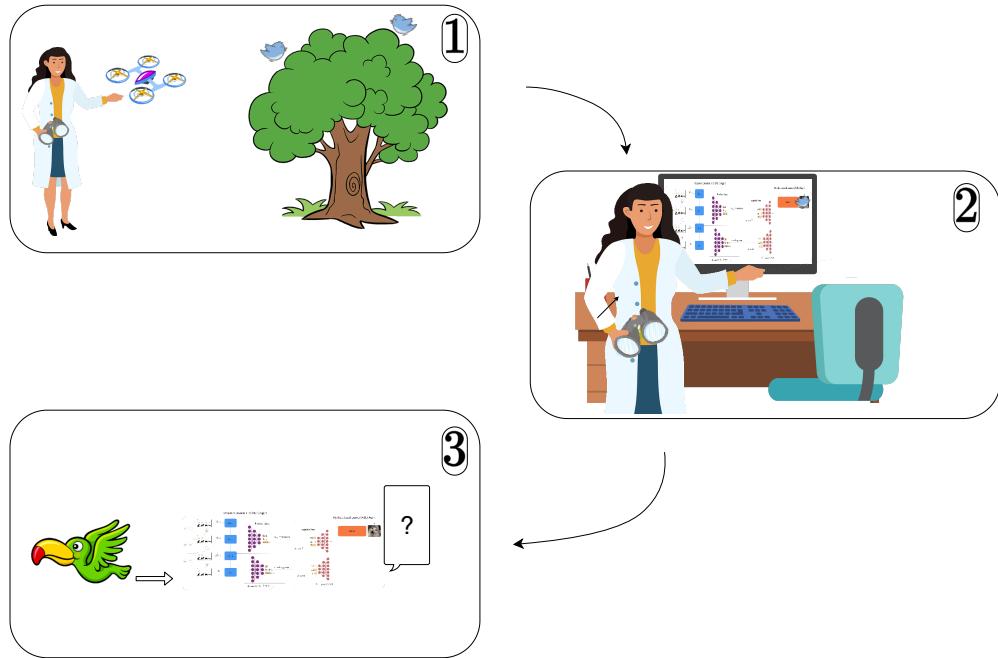


Figure 5.1: After collecting images of birds the ornithologist lets our proposed model classify the vast amount of data. Only in cases of high uncertainty, she is consulted and can classify the image manually.

one epoch was approximately 6 minutes, or in 1 hour, a model can be trained for 10 epochs. Thus for 2000 epochs, we get 200 hours of training time. Considering the specification of our GPUs, we have consumed  $200h \cdot 300\text{ kW} = 60000\text{ kWh}$  of power.

The German federal office for environmental concerns (Umweltbundesamt) estimates the amount of CO<sub>2</sub> necessary to produce one kilowatt-hour to be 401 g [28]. Considering, we required 60000 kWh, this would result in the CO<sub>2</sub> equivalent of  $401\text{ g} \cdot 60000 = 24060\text{ kg}$ .

# Chapter 6

## Conclusion

Extending the work of Alaniz and Akata [2], we propose an explainable model that is aware of and can express its uncertainties. The original model delivers a binary decision tree as an explanation, leveraging human-interpretable attributes, provided as side information. The nodes of the resulting tree correspond to attributes, and the edges to either 'Yes' or 'No' answers regarding the presence of these attributes. In our work, we give the model the ability to either avoid uncertain attributes to be used in the decision tree, or extend it to a ternary decision tree, where '?' is an additional option.

While our proposed extension does not outperform the original aRDTC model, our remRDTC and extRDTC still outperform other explainable methods such as decision trees or dNDF. Moreover, the performance is comparable to that of uninterpretable state of the art methods, such as ResNet.

Our model does not only utilize uncertainty information but also allows introspection. We use this to further investigate the relationship between attribute usage, misclassification and uncertainty. We empirically show that our model (hopefully) reliable yields high uncertainty for unseen attributes. This would allow the model to consult a human user in cases of high uncertainty, thus making it more applicable in real-world scenarios. Finally, we critically reflect the weaknesses of the model and propose possible improvements for future work.

# Bibliography

- [1] Akata, Z., Perronnin, F., Harchaoui, Z., and Schmid, C. (2013). Label-embedding for attribute-based classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 819–826.
- [2] Alaniz, S. and Akata, Z. (2019). Explainable observer-classifier for explainable binary decisions. *arXiv preprint arXiv:1902.01780*.
- [3] Bayes, T. (1763). Lii. an essay towards solving a problem in the doctrine of chances. by the late rev. mr. bayes, frs communicated by mr. price, in a letter to john canton, amfr s. *Philosophical transactions of the Royal Society of London*, (53):370–418.
- [4] Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- [5] Bland, A. R. et al. (2012). Different varieties of uncertainty in human decision-making. *Frontiers in neuroscience*, 6:85.
- [6] Box, G. E. and Tiao, G. C. (2011). *Bayesian inference in statistical analysis*, volume 40. John Wiley & Sons.
- [7] Breiman, L., Friedman, J., Olshen, R., and Stone, C. (1984). Classification and regression trees. statistics/probability series.
- [8] Cao, K., Lazaridou, A., Lanctot, M., Leibo, J. Z., Tuyls, K., and Clark, S. (2018). Emergent communication through negotiation. *arXiv preprint arXiv:1804.03980*.
- [9] d'Amato, A. (1983). Legal uncertainty. *Calif. L. Rev.*, 71:1.
- [10] Damianou, A. and Lawrence, N. (2013). Deep gaussian processes. In *Artificial Intelligence and Statistics*, pages 207–215.
- [11] Das, A., Gervet, T., Romoff, J., Batra, D., Parikh, D., Rabbat, M., and Pineau, J. (2019). Tarmac: Targeted multi-agent communication. In *International Conference on Machine Learning*, pages 1538–1546.
- [12] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.
- [13] Didelot, X., Gardy, J., and Colijn, C. (2014). Bayesian inference of infectious disease transmission from whole-genome sequence data. *Molecular biology and evolution*, 31(7):1869–1879.
- [14] Dikmen, M. and Burns, C. M. (2016). Autonomous driving in the real world:

## Bibliography

- Experiences with tesla autopilot and summon. In *Proceedings of the 8th international conference on automotive user interfaces and interactive vehicular applications*, pages 225–228.
- [15] Farhadi, A., Endres, I., Hoiem, D., and Forsyth, D. (2009). Describing objects by their attributes. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1778–1785. IEEE.
- [16] Foerster, J., Assael, I. A., De Freitas, N., and Whiteson, S. (2016). Learning to communicate with deep multi-agent reinforcement learning. In *Advances in neural information processing systems*, pages 2137–2145.
- [17] Friston, K. J., Glaser, D. E., Henson, R. N., Kiebel, S., Phillips, C., and Ashburner, J. (2002). Classical and bayesian inference in neuroimaging: applications. *Neuroimage*, 16(2):484–512.
- [18] Gal, Y. and Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.
- [19] Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452–459.
- [20] Grossberg, S. (2013). Recurrent neural networks. *Scholarpedia*, 8(2):1888. revision #138057.
- [21] Gunning, D., Stefk, M., Choi, J., Miller, T., Stumpf, S., and Yang, G.-Z. (2019). Xai—explainable artificial intelligence. *Science Robotics*, 4(37).
- [22] Guo, W. (2020). Explainable artificial intelligence for 6g: Improving trust between human and machine. *IEEE Communications Magazine*, 58(6):39–45.
- [23] Guynn, J. (2015). Google photos labeled black people 'gorillas'. *USA Today*.
- [24] Havrylov, S. and Titov, I. (2017). Emergence of language with multi-agent games: Learning to communicate with sequences of symbols. In *Advances in neural information processing systems*, pages 2149–2159.
- [25] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- [26] Hebb, D. (1968). 0.(1949) the organization of behavior.
- [27] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- [28] Icha, P. and Kuhs, G. (2016). Entwicklung der spezifischen kohlendioxid-emissionen des deutschen strommix in den jahren 1990 bis 2015. *Climate Change*, 26:1–27.
- [29] Jiang, J. and Lu, Z. (2018). Learning attentional communication for multi-agent cooperation. In *Advances in neural information processing systems*, pages 7254–7264.
- [30] Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. (1999). An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233.

## Bibliography

- [31] Kearns, M. J. and Mansour, Y. (1998). A fast, bottom-up decision tree pruning algorithm with near-optimal generalization. In *ICML*, volume 98, pages 269–277. Citeseer.
- [32] Kelion, L. (2018). Ces 2018: Lg robot cloi repeatedly fails on stage at its unveil. *BBC*.
- [33] Kendall, A. and Gal, Y. (2017). What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584.
- [34] Koch, K.-R. (2006). *Bayesian inference with geodetic applications*, volume 31. Springer.
- [35] Kotschieder, P., Fiterau, M., Criminisi, A., and Rota Bulo, S. (2015). Deep neural decision forests. In *Proceedings of the IEEE international conference on computer vision*, pages 1467–1475.
- [36] Krzywinski, M. and Altman, N. (2013). Importance of being uncertain.
- [37] Kulkarni, G., Premraj, V., Ordonez, V., Dhar, S., Li, S., Choi, Y., Berg, A. C., and Berg, T. L. (2013). Babytalk: Understanding and generating simple image descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2891–2903.
- [38] Kwon, Y., Won, J.-H., Kim, B. J., and Paik, M. C. (2020). Uncertainty quantification using bayesian neural networks in classification: Application to biomedical image segmentation. *Computational Statistics & Data Analysis*, 142:106816.
- [39] Lampert, C. H., Nickisch, H., and Harmeling, S. (2009). Learning to detect unseen object classes by between-class attribute transfer. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 951–958. IEEE.
- [40] Lazaridou, A., Hermann, K. M., Tuyls, K., and Clark, S. (2018). Emergence of linguistic communication from referential games with symbolic and pixel input. *arXiv preprint arXiv:1804.03984*.
- [41] LeCun, Y., Bengio, Y., et al. (1995). Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995.
- [42] Leibniz, G. W. (1666). *Dissertatio de arte combinatoria, in qua ex arithmeticae fundamentis complicationum ac transpositionum doctrina nouis praeceptis exstruitur... noua etiam Artis meditandis, seu Logicae inuentionis semina sparguntur. Praefixa est synopsis totius tractatus, & additamenti loco demonstratio existentiae Dei, ad mathematicam certitudinem exacta autore Gottfredu Gulielmo Leibniüzio... apud Joh. Simon. Fickium et Joh. Polycarp. Seuboldum in Platea Nicolaea . . .*
- [43] Liu, B. (2013). Toward uncertain finance theory. *Journal of Uncertainty Analysis and Applications*, 1(1):1.
- [44] MacKay, D. J. (1998). Introduction to gaussian processes. *NATO ASI Series F Computer and Systems Sciences*, 168:133–166.
- [45] Marvin, M. and Seymour, A. P. (1969). Perceptrons.
- [46] McCulloch, W. S. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133.

## Bibliography

- [47] Metropolis, N. and Ulam, S. (1949). The monte carlo method. *Journal of the American statistical association*, 44(247):335–341.
- [48] Mingers, J. (1989). An empirical comparison of pruning methods for decision tree induction. *Machine learning*, 4(2):227–243.
- [49] Mullachery, V., Khera, A., and Husain, A. (2018). Bayesian neural networks. *arXiv preprint arXiv:1801.07710*.
- [50] Najafabadi, M. M., Villanustre, F., Khoshgoftaar, T. M., Seliya, N., Wald, R., and Muharemagic, E. (2015). Deep learning applications and challenges in big data analytics. *Journal of Big Data*, 2(1):1.
- [51] NVIDIA, V. (2015). Tesla k80 gpu accelerator board specification.
- [52] NVIDIA, V. (2017). Tesla v100 pcie gpu accelerator data sheet.
- [53] Nwankpa, C., Ijomah, W., Gachagan, A., and Marshall, S. (2018). Activation functions: Comparison of trends in practice and research for deep learning. *arXiv preprint arXiv:1811.03378*.
- [54] O’Connor, K. M. (1994). The alchemical creation of life (takwin) and other concepts of genesis in medieval islam.
- [55] Ordonez, V., Kulkarni, G., and Berg, T. L. (2011). Im2text: Describing images using 1 million captioned photographs. In *Advances in neural information processing systems*, pages 1143–1151.
- [56] Osband, I. (2016). Risk versus uncertainty in deep learning: Bayes, bootstrap and the dangers of dropout. In *NIPS Workshop on Bayesian Deep Learning*, volume 192.
- [57] Palatucci, M., Pomerleau, D., Hinton, G. E., and Mitchell, T. M. (2009). Zero-shot learning with semantic output codes. In *Advances in neural information processing systems*, pages 1410–1418.
- [58] Parmigiani, G. and Parmigiani, G. (2002). *Modeling in medical decision making: a Bayesian approach*. J. Wiley.
- [59] Puiutta, E. and Veith, E. (2020). Explainable reinforcement learning: A survey. *arXiv preprint arXiv:2005.06247*.
- [60] Quinlan, J. (1986). Induction of decision trees. mach. learn.
- [61] Quinlan, J. R. (2014). *C4. 5: programs for machine learning*. Elsevier.
- [62] Raileanu, L. E. and Stoffel, K. (2004). Theoretical comparison between the gini index and information gain criteria. *Annals of Mathematics and Artificial Intelligence*, 41(1):77–93.
- [63] Rodriguez, R. C., Alaniz, S., and Akata, Z. (2019). Modeling conceptual understanding in image reference games. In *Advances in Neural Information Processing Systems*, pages 13155–13165.
- [64] Rosenblatt, F. (1960). Perceptron simulation experiments. *Proceedings of the IRE*, 48(3):301–309.

## Bibliography

- [65] Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *nature*, 323(6088):533–536.
- [66] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252.
- [67] Sarkar, S., Weyde, T., Garcez, A., Slabaugh, G. G., Dragicevic, S., and Percy, C. (2016). Accuracy and interpretability trade-offs in machine learning applied to safer gambling. In *CEUR Workshop Proceedings*, volume 1773. CEUR Workshop Proceedings.
- [68] Schlesinger, A., O’Hara, K. P., and Taylor, A. S. (2018). Let’s talk about race: Identity, chatbots, and ai. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–14.
- [69] Sensoy, M., Kaplan, L., and Kandemir, M. (2018). Evidential deep learning to quantify classification uncertainty. In *Advances in Neural Information Processing Systems*, pages 3179–3189.
- [70] Spiegelhalter, D. and Rice, K. (2009). Bayesian statistics. *Scholarpedia*, 4(8):5230. revision #185711.
- [71] Wagenmakers, E.-J., Love, J., Marsman, M., Jamil, T., Ly, A., Verhagen, J., Selker, R., Gronau, Q. F., Dropmann, D., Boutin, B., et al. (2018). Bayesian inference for psychology. part ii: Example applications with jasp. *Psychonomic bulletin & review*, 25(1):58–76.
- [72] Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. (2011). The Caltech-UCSD Birds-200-2011 Dataset. Technical report.
- [73] Werbos, P. J. (1990). Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560.
- [74] Wu, X., Kumar, V., Quinlan, J. R., Ghosh, J., Yang, Q., Motoda, H., McLachlan, G. J., Ng, A., Liu, B., Philip, S. Y., et al. (2008). Top 10 algorithms in data mining. *Knowledge and information systems*, 14(1):1–37.
- [75] Xian, Y., Lampert, C. H., Schiele, B., and Akata, Z. (2019). Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(9):2251–2265.
- [76] Yadron, D. and Tynan, D. (2016). Tesla driver dies in first fatal crash while using autopilot mode. *The Guardian*.
- [77] Zou, J. and Schiebinger, L. (2018). Ai can be sexist and racist—it’s time to make it fair.

# Acknowledgments

I would like to thank my pet snake, which sleeps all day, and my mother, who is the hardest working person I know. I would like to thank Zeynep and Stephan for the continuous support, their critical questions, and helpful insights.

# **Selbstständigkeitserklärung**

Hiermit versichere ich, dass ich die vorliegende Bachelorarbeit selbstständig und nur mit den angegebenen Hilfsmitteln angefertigt habe und dass alle Stellen, die dem Wortlaut oder dem Sinne nach anderen Werken entnommen sind, durch Angaben von Quellen als Entlehnung kenntlich gemacht worden sind. Diese Bachelorarbeit wurde in gleicher oder ähnlicher Form in keinem anderen Studiengang als Prüfungsleistung vorgelegt.

---

Stefan Wezel (Matrikelnummer 4080589), September 29, 2020