



Technische Universität München

Chair of Media Technology

Prof. Dr.-Ing. Eckehard Steinbach

Master Thesis

Uncertainty-Based Improvement of a Visual
Classification System

Author:	Jianxiang Feng
Matriculation Number:	03679926
Address:	Connollystr. 03/W06 80809 Munich
Advisor:	Prof. Dr. -Ing. Eckehard Steinbach
Supervisor:	Dr. Zoltán-Csaba Márton, Maximilian Durner
Begin:	11.10.2018
End:	25.04.2019

With my signature below, I assert that the work in this thesis has been composed by myself independently and no source materials or aids other than those mentioned in the thesis have been used.

München, April 14, 2019

Place, Date

Signature

This work is licensed under the Creative Commons Attribution 3.0 Germany License. To view a copy of the license, visit <http://creativecommons.org/licenses/by/3.0/de>

Or

Send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California 94105, USA.

München, April 14, 2019

Place, Date

Signature

Abstract

Classifiers based on deep neural networks have achieved tremendous success on different kinds of task. However, at the same time, this kind of classifier could not provide reliable uncertainty estimation about their predictions, which can easily lead to seriously overconfident predictions and sub-optimal decision in the down-stream tasks. From a systematic perspective, this is undesirable and even hazardous in wide range of safety-critical applications such as deployment of robots, medical diagnosis and so on. On the other hand, in object recognition task, some predictions are uncertain because they have similar appearances and thus ambiguous, although their relationship with contextual objects are distinguishable. For example a marker has similar appearance with a toothbrush, when a calculator is recognized nearby, a prediction of marker should have higher probability. This kind of uncertainty from data ambiguity can be considered in order to obtain a more robust classifier.

In this work, regarding the first issue, we investigate applying dropout variational inference and Laplace approximation in Bayesian neural network to obtain reliable model uncertainty estimation on object recognition task. Predictions with low uncertainty can be used to label observed data automatically, which can decrease manual effort on training a more accurate classifier when robot is confronted with objects in real world environment, which has a slight domain gap with objects in training set. When it comes to the second issue, object co-occurrence statistics is introduced as scene contextual information via conditional random field to handle the data ambiguity and improve the exiting results. We have conducted experiments to show that reliable uncertainty estimates can be obtained via dropout inference and even improved with proposed variants of dropout inference and its ensemble. Additionally, those improved uncertainty estimation can be used with semantic statistic to train a more accurate domain specific classifier.

Contents

Contents	ii
1 Introduction	1
1.1 What is uncertainty?	1
1.2 Why do we need uncertainty?	2
1.3 How can we obtain and handle uncertainty in deep learning?	4
1.4 Contributions and Structure	6
2 Bayesian neural network	8
2.1 Introduction	9
2.2 Variational inference	11
2.2.1 Introduction	11
2.2.2 Dropout variational inference	12
2.2.3 Concrete dropout and Multi-Drop	21
2.2.4 Modified network architecture	27
2.3 Laplace approximation	28
2.3.1 Introduction	28
2.3.2 Scalable Laplace approximation for neural network	31
3 Conditional random field	34
3.1 Definition	35
3.2 Learning	37
3.3 Inference	38
4 Experiments and Discussion	40
4.1 Introduction	40
4.1.1 Dataset	40
4.1.2 Uncertainty measures	44
4.1.3 Evaluation metrics	45
4.2 Uncertainty estimation experiments	48
4.2.1 Experiments I	48
4.2.2 Experiments II	50
4.3 Automatic labeling experiments	60

4.3.1	Introduction	60
4.3.2	Experiments I	60
4.3.3	Experiments II	63
4.4	Context-based improvement experiments	63
4.4.1	Experiments settings	63
4.4.2	Results	63
4.4.3	Analysis	63
5	Conclusion	64
5.0.1	Discussions and conclusions	64
5.0.2	Future work	64
A	Appendix	65
A.1	KL condition	65
A.2	Figures	67
A.3	Implementation Details	67
	List of Figures	68
	List of Tables	70
	Bibliography	71

Chapter 1

Introduction

1.1 What is uncertainty?

This question seems a little philosophical and difficult to answer. However, in the context of modeling in machine learning, I find one answer reasonable and illustrative, that is "Uncertainty arises because of limitations in our ability to observe the world, limitations in our ability to model it, and possibly even because of innate nondeterminism." [KFB09]. In spite of innate nondeterminism, uncertainty can be further categorized into two main types :**epistemic uncertainty** and **aleatoric uncertainty** [DKD09][SBD⁺14] [KG17].

The former one, epistemic or model uncertainty illustrates the uncertainty in modeling, which is associated with **imperfect models** of the real world because of insufficient or imperfect knowledge of reality, which corresponds to our ability of modeling, and thus refers to uncertainty caused by a lack of knowledge, i.e., it refers to the epistemic state of the decision maker. Therefore it may be reduced through collecting more knowledge. Examples for this kind of uncertainty can be uncertainty of model parameter, or uncertainty of model structure depending on which kind of model and so on. One simple example to illustrate this kind of uncertainty would be linear regression with limited observed points, where we can use different curves to explain those points and then predict unobserved points with those curves. If we have observed enough points in some intervals, the models or curves that are able to explain the data are more restricted and if there are not enough points, there would be more curves that can explain these points, which indicates high uncertainty. In reality, it's quite normal to have inadequate data, and thus difficult to model the underlying distribution quite well and precisely. Therefore we could explain the data with several or many different models, where the model uncertainty comes in.

The latter one, aleatoric or data uncertainty, accounts for the **inherent randomness** of underlying phenomenon which is expressed as variability in observed data and meanwhile associated with the limitation of our observing ability. It's treated as non-reducible within

our ability to observe the data. Examples for this kind of uncertainty could be noise induced by limited resolution of sensor or ambiguity from content of data which are unrelated to the model we use to represent these data. One example can be classification of images captured by cameras with different resolutions, images with lower resolutions would have higher uncertainty because they are more ambiguous and less illustrative compared with those with higher resolutions.

In other words, epistemic uncertainty refers to the reducible part of the (total) uncertainty, whereas aleatoric uncertainty refers to the non-reducible part. Model uncertainty and aleatoric uncertainty are combined and expressed in the final prediction, which is so called predictive uncertainty.

1.2 Why do we need uncertainty?

Aforementioned model uncertainty represents the belief of model about its predictions and thus bring us more useful information about predictions related with model and observed data instead of just a crispy prediction, whose predictor, is always overconfident when advanced neural network architectures are employed[GPSW17]. As far as I am concerned, this kind of uncertainty plays an important role in two following scenarios.

The first one are those requiring **high safety guarantee** or **optimal decision-making**, in which false predictions can cause hazardous consequences, which is so called AI safety[AOS⁺16]. With reliable model uncertainty estimation, we are able to know when the model is uncertain about its prediction and then people or operator are aware of that and can take corresponding counter measures in order to avoid unnecessary accidents. With more and more wide-spread deep learning algorithms in real life applications, concrete examples for this case include steering control in self-driving car[MGK⁺17], disease diagnosis in medical domain[LAA⁺17], or even tasks in aerospace or other domains requiring higher precision and stronger robustness and so on. Although someone would argue that if the trained model can perform perfectly, that means that it is able to achieve 100% accuracy, then we do not need this kind of uncertainty information. However, as mentioned in [DL91], it's hard or even impossible to have enough training data to define a precise model, which is in coincidence with reality because it's hard to define a task very precisely and get access to enough samples drawn from its real distribution. Uncertainty information would be more valuable and important when algorithms are confronted with multiple tasks[KGC]. Other relevant tasks such as misclassified and out-of-distribution data detection[HG16] and adversarial attacks [KGB16][FCSG17] have also attracted more and more research interests in uncertainty estimation in order to address these challenges.

The second one are scenarios requiring **interaction between algorithms and people or environments**. Model uncertainty can build a bridge between machine learning algorithms and human beings which can construct a more robust and more data-efficient

machine learning system. This scenario involves different kinds of tasks such as active learning[GIG17], reinforcement learning[BCKW15][OBPVR16][GMR], automatically labeling, industrial components inspection and so on. The idea behind that is, with reliable model uncertainty estimation we know which data sample or prediction the model is unfamiliar with. In active learning or reinforcement learning, those unfamiliar data samples can be used to train our model more quickly and efficiently. In tasks involving human robot interaction, we can know when the model requires help from human experts based on uncertainty estimation. For example, in industrial component inspection task, we firstly train a model which can perform quite well on some evaluation datasets and then deploy them in real applications. At this moment, we assume that all predictions made by the model are true all the time. However, there is no guarantee that this model has learned to represent the real distribution of data perfectly and is robust against slight or even large domain transfer. If the distributions of the real world data vary because of unexpected factors such equipments aging or changes of weather condition. The model would fail silently and lead to unexpected accidents. With reliable model uncertainty estimation, this kind of weakness can be mitigated or even eliminated.

On the other hand, aleatoric uncertainty mainly assist in **improving model performance** by quantifying and considering this kind of information by representing noisy level of predictions and making good use of data points with less noise. By taking this information into account, the model can be trained more optimally and yield better performance, some examples of image data and Lidar data can be referred to [KC16][FRD18]. Instead of representing noise, data uncertainty also includes uncertainty from semantic view point, which we can also incorporate into training and yield more robust model. This kind of uncertainty is not quantified explicitly but can be incorporated conceptually with help of specifically designed model and high level statistics cues or features to **disambiguate prediction** directly. One example of this can be modeling contextual information among pixels with conditional random field in semantic segmentation task[KK11][SM⁺12][LSVDHR16]. The idea behind that is simple and conditional random field provides a principle way to achieve this idea. Since uncertainty caused by appearance ambiguity represents strong correlation between different categories. Conditional random field can model correlation or dependency between random variables by taking into account high order potentials. In semantic segmentation task, because of the reason of tractable computation, at most second order potentials are taken into account and can yield better performance compared results with just first order potentials.

1.3 How can we obtain and handle uncertainty in deep learning?

As we know, deep learning is a powerful black box model while it has achieved tremendous success in different tasks. Therefore to obtain and handle uncertainty in deep learning model and data can help us to get more understanding about this black box model and make deep learning based application more robust and safer.

Regarding model uncertainty or quality measure that acts similar to model uncertainty, there are many different ways to obtain this kind of quality measure of prediction. On the one hand, because of intractability of real model uncertainty in deep neural networks, there are some **ad-hoc approaches** that try to obtain such a quality measure on specific tasks which are solved usually with model uncertainty such as misclassified and out-of-distribution(OOD) detection task. [LLS17] combines temperature scaling and input preprocessing which is called adversarial training together and yields a simple but effective out-of-distribution sample detector, but this approach requires additional out-of-distribution dataset for training which is hard to collect in reality. [DT18] modifies the loss function to train a classifier against OOD data, however this approach is hard to generalize in case of slight domain transfer. [LLS17] treats uniform distribution as ground truth distribution of OOD data and show that OOD data can be generated via generative adversarial network(GAN). This approach seems practical, but training GAN effectively does not always work for images of large size because of high dimension of data distribution.

On the other hand, there are some more theoretically sound approaches which employ **Bayesian neural network**[Mac92][Nea12], **bootstrap**[OBPVR16] or **ensemble method**[LPB17] or even ensemble of Bayesian neural network[SG18]. Although Bayesian neural network provides a principal way to obtain model uncertainty by putting randomness on model parameters, it's difficult to infer the exact posterior distribution over model parameters of deep neural network because of its complex architecture and large scale dataset. Therefore approximate inference plays an important role in addressing this issue. The golden standard for approximate inference is Hamiltonian Monte Carlo[Nea12], which requires dealing with whole batch training data and storing the samples of posterior distribution. While the former issue is addressed partially by stochastic gradient Langevin Dynamics[WT11], the latter one attracts many research interests like [BRMW15][WVL⁺18]. However, since it's a Markov Chain Monte Carlo(MCMC) method, to assess the convergence is still non-trivial.

There are some alternatives to MCMC, one popular one is variational inference(VI)[HVC93]. In VI, this inference problem is cast into optimization problem by minimizing the Kullback-Leibler divergence between the approximate posterior distribution and real posterior distribution. Therefore approximate distribution family needs to be chosen and thus has a large impact on the result, which may impede the performance if the chosen distribution family is not flexible enough for the problem. [Gra11]

firstly employed fully factorized Gaussian with a biased gradient estimator on a practical problem. Later [BCKW15] improved the result with the same approximate distribution but different estimator with help of "reparameterization trick" [KW13]. HLA15 also used the same approximate distribution but with expectation propagation [Min01] instead of variational inference. Since the fully factorization assumption may impose an restriction on the flexibility of the approximate distribution between the correlations between weights could not be learned and expressed. By taking this into account, there are many attempts trying to use more expressive approximate distribution such as Bernoulli (treated as mixture of two Gaussian with small variance) [GG16] and Gaussian [KSW15] which capture variances of rows in weight matrix. Another more expressive one is matrix variate Gaussian distribution [LW16] [SCC17] [ZSDG17], which capture both rows and columns correlation in weight matrix. Additionally inspired by the idea of normalizing flow in latent variable models [LW17] applied normalizing flows to auxiliary latent variables to produce more flexible approximate posteriors.

Another alternative is Laplace approximation [Mac92], which put a Gaussian centered at MAP estimate of model parameters with a covariance determined by the inverse of Hessian of the log-likelihood. The main issue of this approach for the large neural network is the large size of Hessian which is quadratic of the number of parameters and inversion operation on it whose complexity is cubic over number of elements of Hessian matrix. Recently [RBB18] working out this problem by approximating the Hessian or the Fisher matrix of log-likelihood with kronecked-factorization which reduces the storage size and complexity of inversion operation. The resulting posterior can also be treated as matrix variate Gaussian.

However, most of those ideas are verified with a relatively simple architecture and on toy or small-scale datasets and need to be modified a lot in order to serve the existing advanced architectures, which is undesirable in our work and many practical applications. Therefore dropout inference [GG16] looks promising considering this restriction and thus is adopted to obtain model uncertainty in our work. On the other hand, scalable Laplace approximation [RBB18] has the potential resolving this issue, because it requires only one MAP point estimate of model parameters, so there is no need to modify the training phase which could save a lot of computation effort and unnecessary modifications for many existing architectures.

When it comes to aleatoric uncertainty, as noise of data itself, it can be modeled by adding another head on top of the network [KG17]. However, the uncertainty caused by appearance ambiguity is not trivial to include in deep neural network training. Therefore another model, **conditional random field** [LMP01], is adopted based on the predictions of discriminative neural network classifier. As is mentioned in previous section, choice of features of second order potentials is key factor in this problem. In this work, because we work on scene object recognition task, therefore the correlation between objects in specific scene are of importance. The correlation can be represented by co-occurrence matrix co-occurrence [LRKT10] [RV09] [GRB08] [RVG⁺07].

1.4 Contributions and Structure

This work can be divided into two main parts, both focusing on uncertainty-based improvement for a deep learning based classifier on object recognition task.

The first part is to employ dropout inference to obtain reliable and calibrated uncertainty from classifier with ResNet50[HZRS16] as backbone. We evaluate the prediction and calibration performance of this probabilistic classifier and then based on uncertainty estimation, we select predictions with low uncertainty/high confidence as automatically labeled data which is used for training a more accurate classifier while the performance of original one declines a lot in case of slight domain transfer which is caused by possibly by light condition changes, appearance changes and even recording equipment changes. This is proof-of-concept idea which tries to proof that improved uncertainty estimation can assist in automatic labeling and reducing manual effort in fine-tuning a new classifier in case of slight domain transfer during deployment of robot. In this part, we consider advanced version of dropout inference and its variants for improving the results both in terms of accuracy and calibration.

The second part is to handle the uncertainty caused by appearance ambiguity. To note that again, this kind of uncertainty is not modeled directly, but the concept is taken into account in choosing model and features. To achieve this goal, we choose conditional random field to model the correlation between objects with similar appearance but different semantics. As key factor of the model, we employ co-occurrence statics of categories as features of second order potentials.

There are several main contributions of this work:

- to show that dropout inference and Laplace approximation can obtain better uncertainty information on multi-class classification problem with 51 classes, which can be further improved by learning dropout rates during training and its ensemble.
- to show that manual effort in data collection in slight domain transfer learning task can be reduced by making use of reliable uncertainty estimation.
- to show that uncertainty caused by appearance ambiguity can be resolved by integrating contextual information via condition random field.

The structure of this work is as follows, in order to give reader a clear theoretical background of models employed in this work, I will briefly review the basic theory and ideas of Bayesian neural network as well as dropout inference and Laplace approximation in chapter 2, where detailed adaptations and modifications of original model are described. In chapter 3, similar to previous chapter, the theory of conditional random field is reviewed firstly and the adaptations are described. In chapter 4, different experiments and results about evaluation of different variants of dropout inference are given and discussed firstly. In the following, experiments and results on scene recognition task resolved by conditional random field with co-occurrence statics as second order features are also showed and discussed.

Then conclusion and possible future work come in the chapter 5.

Chapter 2

Bayesian neural network

As is mentioned in the introduction chapter, we investigate obtaining model uncertainty via Bayesian neural network. Therefore in this chapter, a brief review on theory of Bayesian neural network is given. Since capturing the exact posterior distribution over the parameters of neural network is intractable, we need to have techniques to perform approximate inference. Those techniques could be categorized into two main types, variational inference(VI) and Markov Chain Monte Carlo(MCMC). Brief history and basic idea of these two types of method are introduced. However, since larger dataset and more complex neural network are applied nowadays, traditional approximate inference techniques are difficult to scale to large dataset and complex advanced architectures. Therefore modern approximate inference needs to adopt to obtain model uncertainty. We will introduce and evaluate two modern variational inference methods for neural network, which are dropout inference[GG16] and scalable Laplace approximation[RBB18]. Additionally, we investigate improving uncertainty estimation by learning dropout rate from data via concrete dropout [GHK17], which is introduced following section of dropout inference. Finally, with help of concrete dropout, a new variants of dropout inference is introduced here, which has more flexible approximation distribution family.

2.1 Introduction

Let \mathcal{D} denote an observable dataset consisting of a set of input and output pairs, that is $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\} = \{(x_i, y_i)_{i=1}^N\}$, where $x_i \in \mathbb{R}^D$ and $y_i \in \mathbb{Y}$, \mathbb{Y} is set of labels, and f^ω denote one parametric model which is in our case neural network with ω its parameters which are weights $W_{1:L}$ and biases $b_{1:L}$ for L layers.

In supervised learning, our goal is to learn a probabilistic model of conditional distribution of output given input, $p(\mathbf{y}|\mathbf{x})$ that can explain the underlying data distribution instead of just the observables \mathcal{D} very well based on the observables \mathcal{D} which are samples drawn from the underlying data distribution. Then this learned model can be used to make predictions on unobservable data samples under the same data distribution. In classification where output is label represented by discrete integer, the likelihood function is defined based on parametric model, which is softmax score:

$$p(y = d|\mathbf{x}, \omega) = \frac{\exp(f_d^\omega)}{\sum_{d'} \exp(f_{d'}^\omega)} \quad (2.1)$$

In regression case, the likelihood is Gaussian:

$$p(\mathbf{y}|\mathbf{x}, \omega) = \mathcal{N}(\mathbf{x}; f^\omega(\mathbf{x}), \tau^{-1}\mathbf{I}) \quad (2.2)$$

with model precision τ , which represents the inverse of noise level of the outputs.

As mentioned above, learning means to find model(s) which is parameterized by a set of parameters that can explain the data well, which means the likelihood should be maximized w.r.t. model parameters over the observables. On the other hand, some prior constraints are imposed on the model via prior distribution of the model parameters $p(\omega)$. The probability distribution of model parameters is updated from prior distribution into posterior distribution after observing training dataset \mathcal{D} via Bayes' theorem:

$$p(\omega|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)}{p(\mathbf{Y}|\mathbf{X})} \quad (2.3)$$

where $p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)d\omega$ is so called model evidence or marginal likelihood, whose integration is always intractable.

After obtaining posterior distribution over model parameters, we can make predictions by marginalizing the likelihood of unseen input points such as x^* over model parameters, which leads to predictive distribution over output:

$$p(y^*|x^*, \mathcal{D}) = \int p(y^*|x^*, \omega)p(\omega|\mathcal{D})d\omega \quad (2.4)$$

To point out the difference between normal deterministic neural network and Bayesian neural network can help understanding the mechanism of Bayesian neural network better. In figure 2.1, we use graphical model to express these two kinds of neural network, in which solid point denotes deterministic variable, and circle denotes random variable, while shaded circle denote observed random variable which represents training set. Plate notation denotes N observed data pairs. As we can see in the figure, parameters ω in deterministic is normal variable which is a point estimate done by Maximum a posterior method:

$$\omega^* = \operatorname{argmax}_{\omega} \{p(\mathbf{Y}|\mathbf{X}, \omega)\} \quad (2.5)$$

On the other hand, in Bayesian neural network, the model parameter ω is random variable which obey the distribution parameterized by θ . To note that for explanation of concept here, we do not make any assumption on the function parameterized by θ , which means it can be arbitrary function. This distribution over model parameters can be inferred based on Bayes' rule in equation 2.3. However, for model with large number of parameters, it's hard to calculate the integral tractably. Therefore many approximation methods mentioned in the introduction chapter are applied to solve this problem. It's worth to note that if we choose the approximate distribution as delta function, $q_{\theta}(\omega) = \delta(\omega - \theta)$, then we can recover Bayesian neural network into deterministic neural network. Because the computation of model evidence $p(\mathbf{Y}|\mathbf{X})$ is always intractable in complex model and large dataset, we need to resort to approximate inference, which will briefly be introduced in next section.

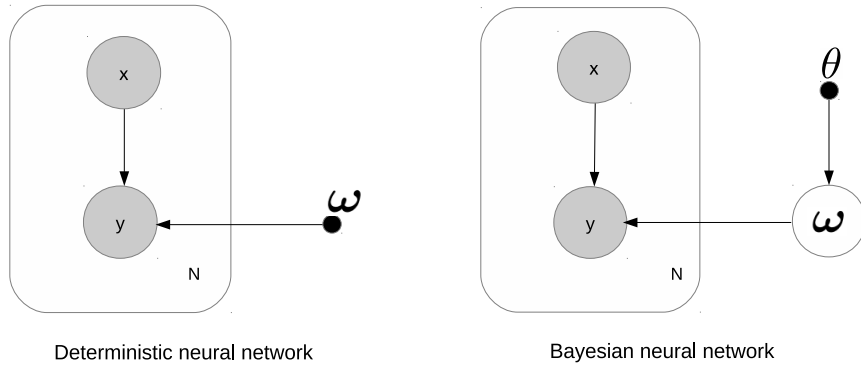


Figure 2.1: Difference between parameter estimation of deterministic neural network and Bayesian neural network.

2.2 Variational inference

2.2.1 Introduction

As is mentioned above, variational inference cast inference into optimization by minimizing the Kullback-Leibler divergence between approximate posterior distribution and the real posterior distribution. However, there is no analytical definition of this KL divergence between the real posterior distribution is unknown. We can derive a lower bound which is also called evidence lower bound (*ELBO*) which bounds the log marginal likelihood with Jensen's inequality. And from that we know that marginal likelihood is the sum of *ELBO* and KL divergence between approximate posterior and real posterior. The derivation is given in the following:

$$\begin{aligned}
 \log(p(\mathbf{Y}|\mathbf{X})) &= \log\left(\int p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})d\boldsymbol{\omega}\right) \\
 &= \log\left(\int q_{\theta}(\boldsymbol{\omega})\frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{q_{\theta}(\boldsymbol{\omega})}d\boldsymbol{\omega}\right) \\
 &\geq \int q_{\theta}(\boldsymbol{\omega})\log\left(\frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{q_{\theta}(\boldsymbol{\omega})}\right)d\boldsymbol{\omega} \\
 &= \mathbb{E}_{q_{\theta}(\boldsymbol{\omega})}[\log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}))] - KL(q_{\theta}(\boldsymbol{\omega})||p(\boldsymbol{\omega})) \\
 &= ELBO
 \end{aligned} \tag{2.6}$$

where $p(\mathbf{Y}|\mathbf{X})$ is the likelihood, $p(\boldsymbol{\omega})$ is the prior distribution over model parameters, $q_{\theta}(\boldsymbol{\omega})$ is the approximate posterior distribution over parameters which is parameterized by θ .

We can get $\log(p(\mathbf{Y})|\mathbf{X})$ by calculating the sum of *ELBO* and KL divergence between approximate posterior $q_{\theta}(\boldsymbol{\omega})$ and real posterior $p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})$:

$$\begin{aligned}
 &ELBO + KL(q_{\theta}(\boldsymbol{\omega})||p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})) \\
 &= \int q_{\theta}(\boldsymbol{\omega})\log\left(\frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{q_{\theta}(\boldsymbol{\omega})}\right)d\boldsymbol{\omega} + \int q_{\theta}(\boldsymbol{\omega})\log\left(\frac{q_{\theta}(\boldsymbol{\omega})}{p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})}\right)d\boldsymbol{\omega} \\
 &= \int q_{\theta}(\boldsymbol{\omega})\log(p(\mathbf{Y}|\mathbf{X}))d\boldsymbol{\omega} \\
 &= \log(p(\mathbf{Y}|\mathbf{X}))
 \end{aligned} \tag{2.7}$$

When we maximize the *ELBO* w.r.t the parameter of approximate posterior θ , which is also called **variational parameter**. We use these two terms interchangeably in the following. It's equivalent to minimizing the KL divergence because the log marginal likelihood is not a function of θ . Then we have a well-defined objective which is the *ELBO*, in which the first term is called expected log likelihood which ensures the model can explain the data well and the second term is called regularization term which ensures the approximate posterior

does not deviate too far from the prior distribution. Now we have cast an inference problem into an optimization problem:

$$\theta^* = \operatorname{argmin}_{\theta} [KL(q_{\theta}(\boldsymbol{\omega}) || p(\boldsymbol{\omega} | \mathbf{X}, \mathbf{Y}))] \quad (2.8)$$

which is equivalent to

$$\theta^* = \operatorname{argmax}_{\theta} [ELBO] \quad (2.9)$$

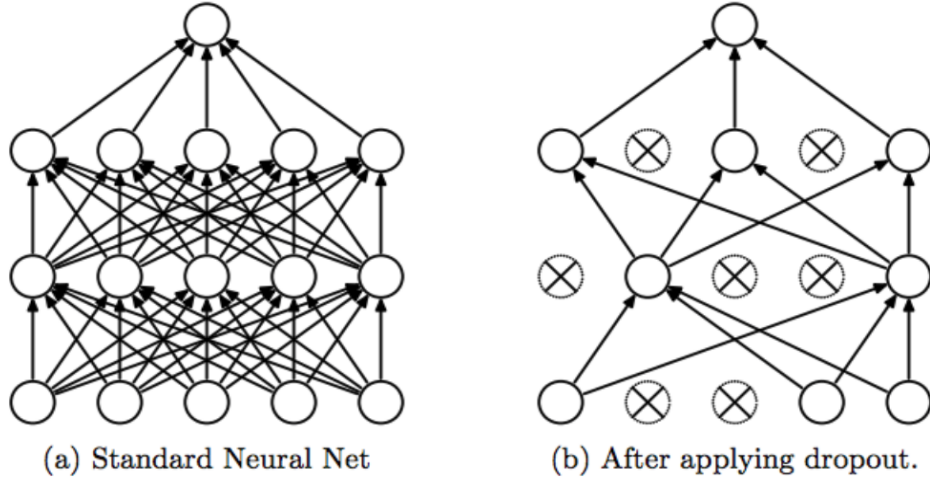
However, there are still some difficulties if we want to solve this optimization with this objective. The first one is to deal with large size dataset, which induces large computation in the expected log likelihood term. [Gra11] shows that this can be solved by data sub-sampling which is also called stochastic optimization. Another one is that we need to obtain the derivatives of $ELBO$ w.r.t. approximate parameter θ . Since the model parameters are samples from approximate distribution, good estimator for the derivatives is required which will be introduced in next subsection.

2.2.2 Dropout variational inference

Dropout

Dropout[SHK⁺14] is originally introduced as regularization approach in training deep neural network which can improve the generalization performance. Although the author said this idea is inspired from human beings sexual reproduction, there are different interpretations trying to explain why it can work such as ensemble perspective and Bayesian perspective. In this subsection, the Bayesian interpretation of dropout is introduced and used for improving the uncertainty estimation of neural network.

The mechanism of dropout is simple, each unit of specific layer is multiplied by a random variable under Bernoulli distribution with $1 - p$ as its parameter, where p is dropout rate. In each iteration of training, dropout is turned on, which means each unit is multiplied by the sample drawn from Bernoulli distribution in forward propagation which is kept in derivatives back propagation during current iteration. In testing, the sampling phase is turned off, only one forward propagation is needed to obtain predictions. Normally, in order to avoid rescaling weights in testing, which is used to keep the output magnitudes in the same scale when dropout is off, rescaling of the output of dropout is always done in training. In figure 2.2, there are two figures about that dropout is on and off.

Figure 2.2: How dropout works[SHK⁺14].

Bayesian interpretation of dropout

As is mentioned in the last subsection, in variational inference, we want to minimize the KL divergence between approximate distribution and the real posterior distribution over the model parameters, which turns out to be maximization of evidence lower bound (*ELBO*). When dropout is interpreted in Bayesian way, the distribution over hidden units is reformulated as distribution over weight matrices. The training objectives of neural network with dropout is proved to be similar as the *ELBO* of Bayesian neural network with Bernoulli distribution factorized over the input dimension of weight matrix. In the following, we will explain this interpretation including following key factors: **approximate distribution**, **training objective**, **marginalization in testing** by using one simple example in classification case in 2.3, in which we define the hidden layer as the first layer and output layer as the second layer and assume that we use L2 regularization and put a prior distribution of fully factorized Gaussian over weights initially, which can be generalized to multi-layer neural network easily.

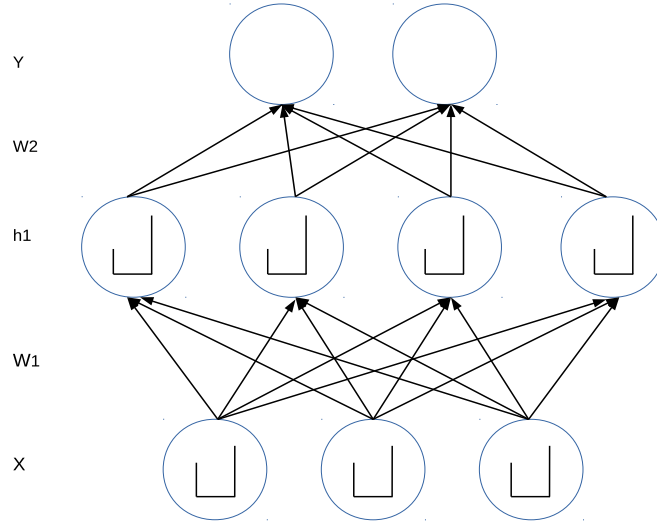


Figure 2.3: A two layer neural network example of dropout inference, a Bernoulli random variable is imposed on each unit of input layer and hidden layer.

Approximate distribution Let's denote $\mathbf{y} \in \mathbb{R}^{m \times D_2}$ as output, $\mathbf{x} \in \mathbb{R}^{m \times D_0}$ as input, $\mathbf{h}_1 \in \mathbb{R}^{D_1}$ as the response of hidden layer, where m represents number of data instances, D_i is dimensionality of i -th layer, where 0-th layer represents input layer and $i \in \{1, \dots, L\}$, $L = 2$ in this example. Further we define $\boldsymbol{\omega} = \{(\mathbf{W}_i)_{i=1}^L\}$ as model parameters, and $\boldsymbol{\epsilon}_i \in \mathbb{R}^{D_{i-1}}$ as the Bernoulli distributed random vector parameterized by $\mathbf{p}_i \in \mathbb{R}^{D_{i-1}}$, for i -th layer. In normal dropout, elements in vector \mathbf{p}_i have same values p_i , which means that there is one keep rate or $(1 - \text{dropout rate})$ for each layer. In the following, we will use \mathbf{p}_i and p_i interchangeably depending on the context if there is no special specification. To note that since weight matrix is treated as random variable, we use $\mathbf{M}_i \in \mathbb{R}^{D_{i-1} \times D_i}$ to denote position of non-zero element in Bernoulli distribution for \mathbf{W}_i . To note that bias $\mathbf{b}_i \in \mathbb{R}^{D_i}$ is absorbed into \mathbf{W}_i by appending a new row at the end of weight matrix and 1 at the end of each data input vector, which is also called homogeneous coordinate. We also assume that approximate weight distribution is factorized over layer, which yields

$$q_{\theta}(\boldsymbol{\omega}) = \prod_{i=1}^L q_{\theta}(\mathbf{W}_i).$$

To start with the formulation of dropout, we model likelihood of output conditioned on

input with softmax scores of neural network in classification case:

$$\begin{aligned}
p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) &= \sigma((\mathbf{h}_1 \odot \boldsymbol{\epsilon}_2) \mathbf{M}_2) \\
&= \sigma(\mathbf{h}_1 (\text{diag}(\boldsymbol{\epsilon}_2) \mathbf{M}_2)) \\
&= \sigma(\mathbf{h}_1 \mathbf{W}_2) \\
&= \sigma(\phi(\mathbf{x} (\text{diag}(\boldsymbol{\epsilon}_1) \mathbf{M}_1) + \mathbf{b}_1) \mathbf{W}_2) \\
&= \sigma(\phi(\mathbf{x} \mathbf{W}_1) \mathbf{W}_2)
\end{aligned} \tag{2.10}$$

where \odot is Hadamard product (element-wise product), $\sigma(a_j) = \frac{\exp(a_j)}{\sum_k \exp(a^k)}$ is softmax function, $\phi(\cdot)$ is element-wise non-linear activation function such as rectified unit function.

From the equation above, we have

$$\begin{aligned}
\mathbf{W}_i &= g(\mathbf{M}_i, \boldsymbol{\epsilon}_i) = \text{diag}(\boldsymbol{\epsilon}_i) \mathbf{M}_i \\
\text{with } \boldsymbol{\epsilon}_i &\sim p(\boldsymbol{\epsilon}_i) = \text{Bernoulli}(\mathbf{p}_i)
\end{aligned}$$

which means weight matrix \mathbf{W}_i is a random variable whose probability density function is parameterized by \mathbf{p}_i and \mathbf{M}_i , which are denoted by $\theta = \{(\mathbf{M}_i, \mathbf{p}_i)_{i=1}^L\}$ in equation 2.8, where $i = 1, \dots, L$ denotes i -th layer of the network, and $L = 2$ in this example. The expression of approximate posterior distribution is not obvious, but we can define the its form as

$$\begin{aligned}
q_\theta(\boldsymbol{\omega}) &= \prod_{i=1}^L q_\theta(\mathbf{W}_i) \\
&= \prod_{i=1}^L \int q_\theta(\mathbf{W}_i | \boldsymbol{\epsilon}_i) p(\boldsymbol{\epsilon}_i) d\boldsymbol{\epsilon}_i
\end{aligned} \tag{2.11}$$

with

$$\begin{aligned}
q_\theta(\mathbf{W}_i | \boldsymbol{\epsilon}) &= \delta(\mathbf{W}_i - g(\mathbf{M}_i, \boldsymbol{\epsilon}_i)) \\
&= \delta(\mathbf{W}_i - \text{diag}(\boldsymbol{\epsilon}_i) \mathbf{M}_i)
\end{aligned} \tag{2.12}$$

As we can see, the approximate posterior distribution over the parameter matrix puts a same Bernoulli distribution over the input dimension of parameter matrix, which is the row dimension in this example. Meanwhile each element of the same row is multiplied with same realization of the random variable but with different non-zero position, which is corresponding to different expectation of each row element and thus induces correlations between row elements. To make this definition more clear, based on the computation of

expectation and variance of Bernoulli distribution, we can write down the first and second moment of the approximate distributed random variables in the following:

$$\mathbb{E}_{q_\theta}(\mathbf{W}_i) = \mathbf{M}_i \odot \mathbf{P}_i \quad (2.13)$$

where $\mathbf{P}_i = [\mathbf{p}_i, \dots, \mathbf{p}_i] \in \mathbb{R}^{D_{i-1} \times D_i}$.

Covariance matrix of parameter matrix is:

$$[Cov_{q_\theta}(vec(\mathbf{W}_i))]_{jk} = \mathbb{1}[l = m] m_{lq}^i * m_{mn}^i * p_i * (1 - p_i) \quad (2.14)$$

where

$$j = l * D_{i-1} + q$$

,

$$k = m * D_{i-1} + n$$

, which is the linear mapping between element index before and after matrix vectorization. m_{lq}^i is the element of l -th row and q -th column in matrix \mathbf{M}_i , which applies to m_{mn}^i as well. $\mathbb{1}[i = j]$ denotes indicator function, which is equal to 1 only when i is equal to j and otherwise 0. Because $vec(\cdot)$ operation converts matrix $\mathbf{W}_i \in \mathbb{R}^{D_{i-1} \times D_i}$ into column vector $vec(\mathbf{W}_i) \in \mathbb{R}^{D_{i-1} D_i \times 1}$ by stacking the columns of matrix on top of one another. Then it's easy to see that covariance matrix $Cov_{q_\theta}(vec(\mathbf{W}_i)) \in \mathbb{R}^{D_{i-1} D_i \times D_{i-1} D_i}$.

From equation 2.14, we can see that the entire covariance matrix is consisting of $D_i * D_i$ diagonal sub-matrices whose dimensionality are $D_{i-1} \times D_{i-1}$. When observing covariance matrix of parameter matrix \mathbf{W}_i w.r.t. approximate posterior distribution $q_\theta(\mathbf{W}_i)$, we know that samples of row vectors in weight matrix are drawn independently and thus covariance between rows are zero. On the other hand, samples for different weights in the same row are drawn at the same time because they are multiplied by the same realization of ϵ_i , from which covariances between weights within the same row are induced. Therefore, by fitting this approximate distribution to the real posterior distribution weights within the same row can be learned. This means that the approximate distribution family have more flexibility to approximate the real posterior distribution when compared with other common approximate posterior distribution family that assumes distribution for each parameter is independent such as fully factorized Gaussian.

Training objective Up to now, we have only analyzed the approximate distribution of dropout inference. As is mentioned in introduction section of this chapter, we know that we perform optimization of *ELBO* w.r.t. the approximate distribution parameters, in order to obtain a good approximation to the true posterior distribution over parameters, which we can use in testing to obtain more reliable uncertainty estimation. In the following,

we will show that training a neural network with dropout is equivalent to optimizing the *ELBO* w.r.t. variational parameters.

At first, let's define the training objective of neural network with dropout, which is cross entropy between predictive distribution and target distribution plus L2 regularization:

$$L_{dropout} = \sum_{i=1}^N \left[-\log(p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\omega})) \right] + \lambda \left(\sum_{i=1}^L \|\mathbf{W}_i\|^2 \right) \quad (2.15)$$

where N represents the size of entire dataset, λ is L2 regularization coefficient. We want to maximize the likelihood over the entire dataset w.r.t. the model parameter $\boldsymbol{\omega}$, which is known as maximum likelihood estimation. Equipped with L2 regularization and Gaussian prior over parameters, we can obtain a max-a-posterior estimation by minimizing equation 2.15.

Normally we use gradient descent to tune our parameters in training, which requires first derivative of objective w.r.t. model parameters $\boldsymbol{\omega}$. Since nowadays large size of dataset is ubiquitous, which means N in equation 2.15 is too large, we could not obtain exact gradient of entire batch with efficient computation. Therefore we use data of mini-batch to estimate gradient, which is so called stochastic gradient descent(SGD). Apart from making computation tractable, noise of gradient estimation in each mini-batch is helpful for optimization procedure to escape the poor local minimum. The expression of gradients required in each iteration of SGD is given in the following:

$$\frac{\partial L_{dropout}}{\partial \boldsymbol{\omega}} = \frac{N}{K} \sum_{i \in S} \left[-\frac{\partial \log(p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\omega}))}{\partial \boldsymbol{\omega}} \right] + \lambda \frac{\partial (\sum_{i=1}^L \|\mathbf{W}_i\|^2)}{\partial \boldsymbol{\omega}} \quad (2.16)$$

where S is one random subset of entire dataset, and $|S| = K$.

On the other hand, let's have a look at *ELBO* in equation 2.6, if we want to maximize *ELBO*, which is equivalent to minimize negative *ELBO* w.r.t. the variational parameter θ , with SGD. The gradients are computed with the following expression:

$$\frac{\partial (-ELBO)}{\partial \theta} = \frac{N}{K} \sum_{i \in S} \left[-\frac{\partial \mathbb{E}_{q_\theta(\boldsymbol{\omega})} [\log(p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\omega}))]}{\partial \theta} \right] + \frac{\partial KL(q_\theta(\boldsymbol{\omega})||p(\boldsymbol{\omega}))}{\partial \theta} \quad (2.17)$$

In order to calculate two terms in equation 2.17, we need to introduce one technique called re-parameterization trick[KW13] for the first term and one condition called KL condition[Gal] for the second term in the following.

Re-parameterization trick: when we take a close look at the first term in equation 2.17, we know that we need to compute the gradients of expectation w.r.t. the parameters of distribution to which this expectation is subject. That means, we need to estimate the gradients of those parameters because our objective is generated from these parameters randomly. Fortunately, there are different approaches to estimate the gradients of this kind of parameters such as score function or REINFORCE estimator [Wil92], re-parameterization trick [KW13] and so on. As is stated in [KW13], re-parameterization trick has lower variance than score function estimator and is also unbiased. Let's have a quick recap of this technique and see it's already a built-in part in neural network with dropout.

To identify the problem, we can write a more general form of calculus we want to solve in the following:

$$I(\theta) = \frac{\partial}{\partial \theta} \mathbb{E}_{p_\theta(x)}[f(x)] = \frac{\partial}{\partial \theta} \int f(x) p_\theta(x) dx \quad (2.18)$$

Assume that $p_\theta(x)$ can be re-parameterized as $p(\epsilon)$ which is a parameter-free distribution such that random variable x can be generated from a deterministic differentiable function with θ and ϵ as arguments, that is

$$x = g(\theta, \epsilon) \text{ with } \epsilon \sim p(\epsilon)$$

Then we can derive the estimator of the gradients w.r.t. distribution parameters θ with $p(x|\epsilon) = \delta(x - g(\theta, \epsilon))$:

$$\begin{aligned} I'(\theta) &= \frac{\partial}{\partial \theta} \int f(x) p_\theta(x) dx \\ &= \frac{\partial}{\partial \theta} \int f(x) \left(\int p_\theta(x|\epsilon) p(\epsilon) d\epsilon \right) dx \\ &= \frac{\partial}{\partial \theta} \int \left(\int f(x) p_\theta(x|\epsilon) dx \right) p(\epsilon) d\epsilon \\ &= \frac{\partial}{\partial \theta} \int \left(\int f(x) \delta(x - g(\theta, \epsilon)) dx \right) p(\epsilon) d\epsilon \\ &= \frac{\partial}{\partial \theta} \int f(g(\epsilon, \theta)) p(\epsilon) d\epsilon \\ &= \int \frac{\partial}{\partial \theta} f(g(\epsilon, \theta)) p(\epsilon) d\epsilon \\ &= \int f'(g(\epsilon, \theta)) \frac{\partial}{\partial \theta} g(\theta, \epsilon) p(\epsilon) d\epsilon \\ &= \mathbb{E}_{p(\epsilon)} \left[f'(g(\epsilon, \theta)) \frac{\partial}{\partial \theta} g(\theta, \epsilon) \right] \end{aligned} \quad (2.19)$$

From practical point of view, the expectation of last line in expression 2.19 can be approximated with Monte Carlo integration. In dropout inference, we know that if we fix

the dropout rate to $1 - p_i$ in equation 2.12 during training. Then ϵ_i in equation 2.12 is a parameter free random variable and $g(\cdot)$ is a differentiable function w.r.t. its input argument \mathbf{M}_i . Now the variational parameter θ only contains $\{(\mathbf{M}_i)_{i=1}^L\}$, which are exactly the weights to be learned in training neural network with dropout. As a consequence, if we estimate the gradient of $ELBO$ w.r.t. $\{(\mathbf{M}_i)_{i=1}^L\}$, which is the first term in equation 2.17, it's equivalent to calculate the gradient of dropout loss w.r.t. model parameters ω which is the first term in equation 2.16.

KL condition: The second term in equation 2.17 is proved to be equivalent to the second term in equation 2.16 for a large enough number of hidden units when we specify the model prior to be a product of independent Gaussian distributions over each weight with prior length scale l , that is:

$$p(\omega) = \prod_{i=1}^L (p(\mathbf{W}_i))$$

with

$$p(\text{vec}(\mathbf{W}_i)) = \mathcal{N}(0, l^{-2} \mathbf{I}_{D_{i-1}D_i})$$

To establish this condition, we need to make a approximation of the approximate posterior distribution $q_\theta(\omega)$ in equation 2.11, where we approximate $q_\theta(\mathbf{W}_i|\epsilon_i)$ in equation 2.12 as a narrow Gaussian with a very small standard deviation. As we know, $q_\theta(\mathbf{W}_i)$ factorizes over each row of the weight matrix. Then that means $q_\theta(\omega)$ is a mixture of two Gaussians with small standard deviations, and one component fixed at zero:

$$\begin{aligned} q_\theta(\omega) &= \prod_{i=1}^L q_\theta(\mathbf{W}_i) \\ &= \prod_{i=1}^L \prod_{j=1}^{D_{i-1}} q_\theta(\mathbf{w}_{i,j}) \\ &= \prod_{i=1}^L \prod_{j=1}^{D_{i-1}} [p_i \mathcal{N}(\mathbf{m}_{i,j}, \sigma^2 \mathbf{I}_{D_i}) + (1 - p_i) \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{D_i})] \end{aligned} \tag{2.20}$$

where $\mathbf{w}_{i,j} \in \mathbb{R}^{D_i}$ is the j -th row of weight matrix \mathbf{W}_i and p_i is the parameter of Bernoulli distributed random variable of i -th layer. With this, the KL divergence between approximate posterior and prior is KL divergence between mixture of Gaussian and a single Gaussian. In order to keep this report as self-contained as possible, we attach the derivation of KL divergence between mixture of Gaussian and single Gaussian in appendix A.1.

Then we have KL condition in the following:

$$KL(q_\theta(\boldsymbol{\omega})||p(\boldsymbol{\omega})) \approx \sum_{i=1}^L \sum_{j=1}^{D_{i-1}} \left[\frac{p_i}{2} (l^2 \mathbf{m}_{i,j}^T \mathbf{m}_{i,j} + D_i(\sigma^2 - \log(\sigma^2) - 2\log l - 1) - \mathcal{H}(\mathbf{p}_i)) \right] \quad (2.21)$$

with

$$\mathcal{H}(\mathbf{p}_i) = D_{i-1}(-p_i \log p_i - (1 - p_i) \log(1 - p_i))$$

for large enough D_i . If we **fix dropout rate** in training and compute the gradients of KL divergence w.r.t. variational parameters $\theta = \{(\mathbf{M}_i)_{i=1}^L\}$, which is now exactly the same as the model parameters $\boldsymbol{\omega} = \{(\mathbf{W}_i)_{i=1}^L\}$. Then the term entropy of dropout rate can be ignored. We can see that, if we choose $\lambda = \frac{l^2 p_i}{2}$, then it's equivalent to the gradients of regularization term in dropout loss function(cf. equation 2.16):

$$\frac{\partial KL(q_\theta(\boldsymbol{\omega})||p(\boldsymbol{\omega}))}{\partial \theta} \approx \frac{\partial \sum_{i=1}^L \lambda \|\mathbf{M}_i\|^2}{\partial \theta} \quad (2.22)$$

With aforementioned re-parameterization trick and KL condition, we know that the computation of gradients of objective function w.r.t. model parameters in equation 2.16 is equivalent to that of negative *ELBO* w.r.t. variational parameter in equation 2.17. That means, if we train a neural network with dropout, it's equivalent to train a Bayesian neural network with approximate posterior distribution over model parameters defined in equation 2.11.

Marginalization in testing After we have obtain the parameters θ of approximate posterior distribution over model parameters $q_\theta(\boldsymbol{\omega})$, we can marginalize all possible parameters over approximate posterior distribution to get final predictive distribution of output, similar to equation 2.4 but now with approximate posterior distribution over model parameter, which has an expression in equation 2.23. Because the integration is hard to evaluate analytically. In practice, we always use Monte Carlo integration:

$$\begin{aligned} p(\mathbf{y}^*|\mathbf{x}^*, \mathcal{D}) &= \int p(\mathbf{y}^*|\mathbf{x}^*, \boldsymbol{\omega}) p(\boldsymbol{\omega}|\mathcal{D}) d\boldsymbol{\omega} \\ &\approx \int p(\mathbf{y}^*|\mathbf{x}^*, \boldsymbol{\omega}) q_\theta(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &\approx \frac{1}{K} \sum_{i=1}^K p(\mathbf{y}^*|\mathbf{x}^*, \hat{\boldsymbol{\omega}}_i) \end{aligned} \quad (2.23)$$

where $\boldsymbol{\omega} \sim q_\theta(\boldsymbol{\omega})$ and $\hat{\boldsymbol{\omega}}_i$ is one of K realizations of $\boldsymbol{\omega}$, which is equivalent to turning on dropout in testing time. This is also called **MC-dropout** in the literatures.

2.2.3 Concrete dropout and Multi-Drop

Concrete dropout Based on the aforementioned description of dropout inference, we can see that if we fix dropout rate in training, then parameters of approximate distribution θ only contains $\{(\mathbf{M}_i)_{i=1}^L\}$, which is equivalent to $\omega = \{(\mathbf{M}_i)_{i=1}^L\}$ in neural network with dropout. Therefore optimising any neural network with dropout is equivalent to a form of approximate inference in a probabilistic interpretation of the model. This means that the optimal weights found through the optimization of a neural network with dropout are the same as the optimal variational parameters in a Bayesian neural network with the same structure.

However, fixing dropout rate in train is not a trivial task for several reasons. Firstly, as is shown in [SHK⁺14], different dropout rates have different impact on model capacity and thus model performance. To choose an optimal dropout rate manually requires repeating tedious experiments and thus waste of time and computation effort. Secondly, if we want our model not only to achieve satisfied performance but also to possess reliable uncertainty estimation, the dropout rate matters a lot because it belongs to variational parameter θ and further influences the flexibility of approximate distribution family from the perspective of Bayesian interpretation.

Accordingly, one direct counter measure is to learn dropout rate from the data [GHK17]. One major difficulty to learn dropout rate from the data in gradient-based optimization is that it's not trivial to estimate gradients of expectation w.r.t. parameters of a discrete distribution. Before when the approximate distribution is continuous, we estimate this gradient with help of re-parameterization trick. As introduced in the last section, re-parameterization trick requires re-parameterizing samples of the distribution into a differential function with variational parameters and a parameter-free random variable as input arguments. For continuous distribution, this function can be found easily if they have tractable inverse cumulative density function or functional form like Gaussian [KW13]. For most of discrete distributions such as Bernoulli distribution or categorical one, they lack useful re-parameterizations due to the discontinuous nature of discrete states [MMT16].

Confronted with this issue, [JGP16] and [MMT16] come up with one approach where "max" operation in Gumbel-Max trick is replaced by "softmax" function, which can yields a practical re-parameterization for discrete random variable. With this, gradients w.r.t. parameters of discrete approximate distribution can be computed and used in gradient-based optimization. In this subsection, a quick recap of this approach is given in the following.

Re-parameterization of Bernoulli distribution Firstly, Gumbel-Max trick [MTM14] is introduced in figure 2.2.3, which is used for drawing samples from a discrete distribution which is parameterized by set of unnormalized probability $\{\alpha_i\}_{i=1}^n$ via inverse cumulative distribution function of Gumbel distribution, where $\alpha_i \in \mathbb{R}_{>0}$ and n denotes the number of

class. Assuming that we use one-hot encoding vector for the class representation, that is $\mathbf{d} \in \{0, 1\}^n$ and $\sum_{i=1}^n d_i = 1$. The Gumbel-Max trick proceeds as follows(cf. figure 2.2.3):

- sample $G_i \sim \text{Gumbel}(0, 1) = -\log(-\log(\text{Uniform}(0, 1)))$, for $i = 1, \dots, n$
- compute $x_i = \log(\alpha_i) + G_i$, for $i = 1, \dots, n$
- set $d_k = 1$, where $k = \text{argmax}_i \{x_i\}_{i=1, \dots, n}$, and $d_i = 0$ for $i \neq k$

Then we obtain one sample from this discrete distribution and the probability for specific class.

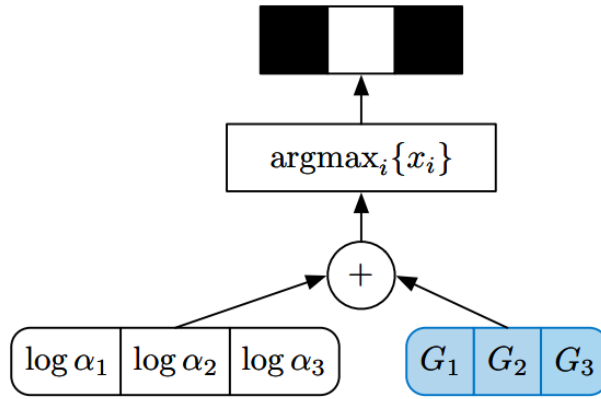


Figure 2.4: Example of Gumbel-max trick for drawing samples from a discrete distribution whose random variable has 3 states and $\{\alpha_i\}_{i=1,2,3}$ as class parameters representing the possibility of occurrence of that class. $\{G_i\}_{i=1,2,3}$ are i.i.d Gumbel(0, 1) [MMT16].

As observed in Gumbel-Max trick, the sampling step is re-paramterized by one function with distribution parameters and parameter-free random variable as arguments. This function is component-wise addition of two input arguments followed by *argmax* operation. However, *argmax* operation is not differential w.r.t. distribution parameter α_i .

With a little abuse of notation, we denote temperature by λ here. Fortunately, the *argmax* operation can be approximated by a continuous function, *softmax* function with λ as temperature parameter, which is differential w.r.t. distribution parameters. With this replacement, we can obtain a continuous approximation to one-hot encoding vector \mathbf{d} on the simplex $\Delta^{n-1} = \{\mathbf{x} \in \mathbb{R}^n | x_k \in [0, 1], \sum_{i=1}^n x_i = 1\}$:

$$x_k = \frac{\exp((\log \alpha_k + G_k)/\lambda)}{\sum_{i=1}^n \exp((\log \alpha_i + G_i)/\lambda)} \quad (2.24)$$

The sampling steps are similar to Gumbel-Max trick, but with smoothed *softmax* operation instead of *argmax*(cf. figure 2.5). When $\lambda \rightarrow 0$, the *softmax* function is recovered

to *argmax* operation, when $\lambda \rightarrow \infty$, this operation generates uniform vector instead of one-hot encoding vector. In practice, we set temperature as 0.1.

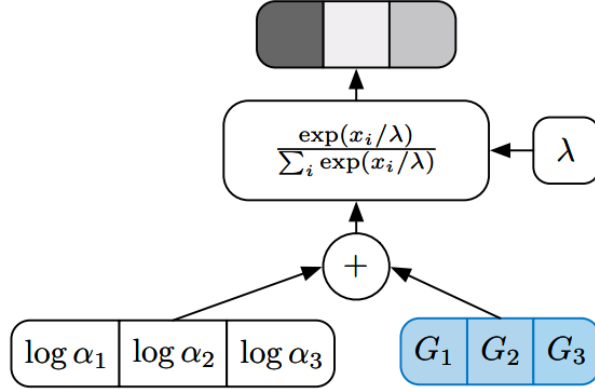


Figure 2.5: Example of continuous approximation to Gumbel-max trick for drawing samples from a discrete distribution whose random variable has 3 states and $\{\alpha_i\}_{i=1,2,3}$ as class parameters representing the possibility of occurrence of that class. $\{G_i\}_{i=1,2,3}$ are i.i.d Gumbel(0, 1)[MMT16].

When it comes to Bernoulli discrete distribution in our case, it becomes simpler because there are only two states in this distribution and samples live in two dimensional simplex. On the other hand, the difference of two Gumbels distributed random variable is similar to a logistic distributed random variable. With the probability of state 1 which can be expressed by:

$$\begin{aligned}\mathbb{P}(d_1 = 1) &= \mathbb{P}(\log \alpha_1 + G_1 > \log \alpha_2 + G_2) \\ &= \mathbb{P}(\log \alpha_1 - \log \alpha_2 + G_1 - G_2 > 0)\end{aligned}$$

where $G_1 - G_2 = \log(\text{Uniform}(0, 1)) - \log(1 - \text{Uniform}(0, 1))$, which is the inverse cumulative density function of Logistic(0,1). Then we can infer the re-parameterization of continuous approximation of sample $x \in [0, 1]$ from Bernoulli distributed random variable with p as parameter as follows:

$$x = \text{sigmoid}\left(\frac{1}{\lambda}(\log p - \log(1 - p) + u - \log(1 - u))\right) \quad (2.25)$$

where $u \sim \text{Uniform}(0, 1)$.

In order to make explanation more clear and validate the continuous approximate re-parameterization, there is one plot (figure 2.6) showing the relationship between average values of 100 samples drawn from this continuous approximation of Bernoulli distribution with different probability parameter from 0 to 1 with step of 0.1. We can see that the

empirical average of samples drawn from equation 2.25 is close to expectation of Bernoulli distribution of corresponding parameter. Additionally, we show one single sample in the figure. As we can see, most of samples are lying on the boundary of range $[0, 1]$ and only few of them lie in the interior.

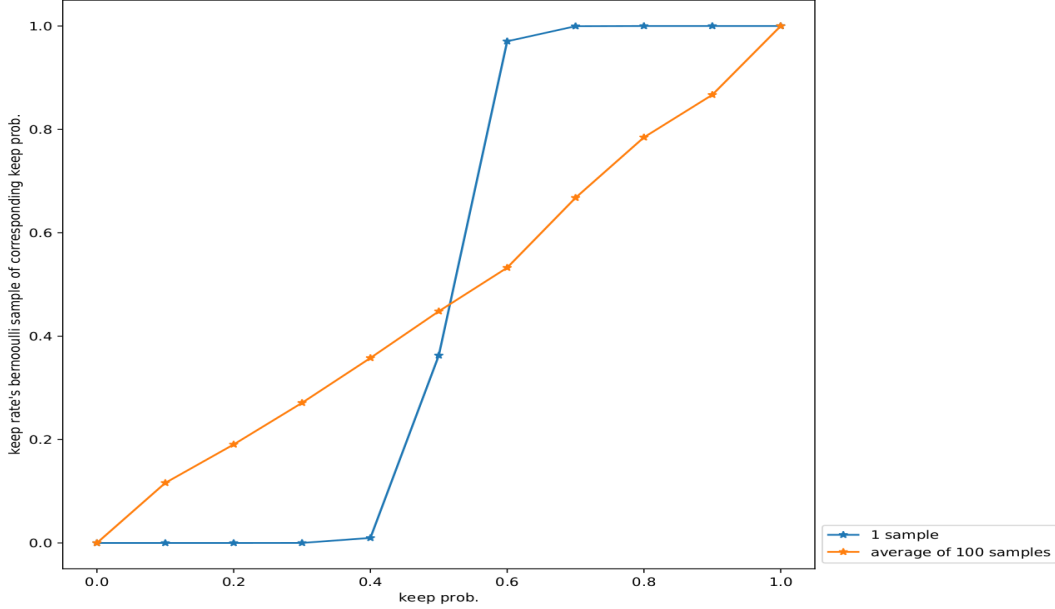


Figure 2.6: One sample and average value of 100 samples drawn from continuous approximation of Bernoulli distribution with parameter $p = [0.1, 0.2, \dots, 1.0]$ and temperature $\lambda = 0.1$.

Dropout regularization Because keep rate of dropout is being optimized, the parameters of approximate distribution θ contains both $\{(\mathbf{M}_i)_{i=1}^L\}$ and $\{(\mathbf{p}_i)_{i=1}^L\}$. In order to compute gradients of ELBO in equation 2.17, we employ categorical re-parameterization to estimate the gradients w.r.t. Bernoulli distribution parameter \mathbf{p}_i in the first term, for the second term, we could not ignore the term with \mathbf{p}_i in KL condition 2.21, which is the entropy term. Consequently, unlike equation 2.22, the KL divergence term should be:

$$\begin{aligned} \frac{\partial KL(q_\theta(\omega) || p(\omega))}{\partial \theta} &\approx \frac{\partial \sum_{i=1}^L \lambda \|\mathbf{M}_i\|^2 - \beta \mathcal{H}(\mathbf{p}_i)}{\partial \theta} \\ &= \frac{\partial}{\partial \theta} \left(\sum_{i=1}^L \lambda \|\mathbf{M}_i\|^2 - \beta D_{i-1}(-p_i \log p_i - (1-p_i) \log(1-p_i)) \right) \end{aligned} \quad (2.26)$$

where λ and β are coefficients for L2 regularization and dropout regularization, respectively. They are hyper-parameters which can be tuned based on performance on validation set.

From the equation above, we can see that this term maximize the entropy of Bernoulli distribution, which means this term pushes \mathbf{p}_i to 0.5. The coefficient of the dropout regularisation term means that large models will push the dropout probability towards 0.5 much more than smaller models, but as the amount of data N increases the dropout probability will be pushed towards 1 because of the expected log likelihood in the first term. One of reasons behind could be pushing \mathbf{p}_i to 0.5 would decrease the capacity of the model which will decrease the expected log-likelihood.

Multi-Drop From figure 2.2 and expression of approximate distribution in equation 2.11, we choose only one probability of Bernoulli distributed random variable for each layer, therefore random vector $\mathbf{p}_i = [p_i]^{D_{i-1}}$ for i -th layer, which stacks same value into a vector. While the dropout regularization term pushes the probability of Bernoulli to 0.5 to maximize its entropy, the expected likelihood term tries to increase the probability because decreasing probability will lead to a different model with lower capacity and thus low performance. An equilibrium state between them should be achieved in training. With concrete dropout introduced above, we could extend dropout for each hidden units instead of each layer(cf. figure 2.8), which means random vector $\mathbf{p}_i = [p_i^k]_{k=1}^{D_{i-1}}$. While the first term in gradients computation stays the same, the second term should be modified to:

$$\begin{aligned} \frac{\partial KL(q_\theta(\boldsymbol{\omega})||p(\boldsymbol{\omega}))}{\partial \theta} &\approx \frac{\partial \sum_{i=1}^L \lambda ||\mathbf{M}_i||^2 - \beta \mathcal{H}(\mathbf{p}_i)}{\partial \theta} \\ &= \frac{\partial}{\partial \theta} \left(\sum_{i=1}^L \lambda ||\mathbf{M}_i||^2 - \beta \sum_{k=1}^{D_{i-1}} (-p_i^k \log p_i^k - (1 - p_i^k) \log(1 - p_i^k)) \right) \end{aligned} \quad (2.27)$$

The reasons behind multi-drop are as following:

- to increase flexibility in tuning variational parameters. The tunability of parameter of Bernoulli random variable in the likelihood term is low because there is only single parameter controlling the entire layer. As is observed in the experiments(cf.figure 2.2.3), these parameters are always increased for each layer. The reason for this is probably that reducing it would lead to low capacity and thus low likelihood.
- the solution space of concrete dropout should be a subset of the solution space of multi-drop if all of them are reachable. Because if it's optimal to assign same probability for each hidden units, this can be recovered in training with multi-drop. Otherwise, other optimal solutions of assigning different probabilities to different hidden units could be considered.
- last but not least, multi-drop can extend the flexibility and diversity of the dropout approximate distribution family by adding more parameters. Hence the truth posterior can be approximated by the approximate posterior better.

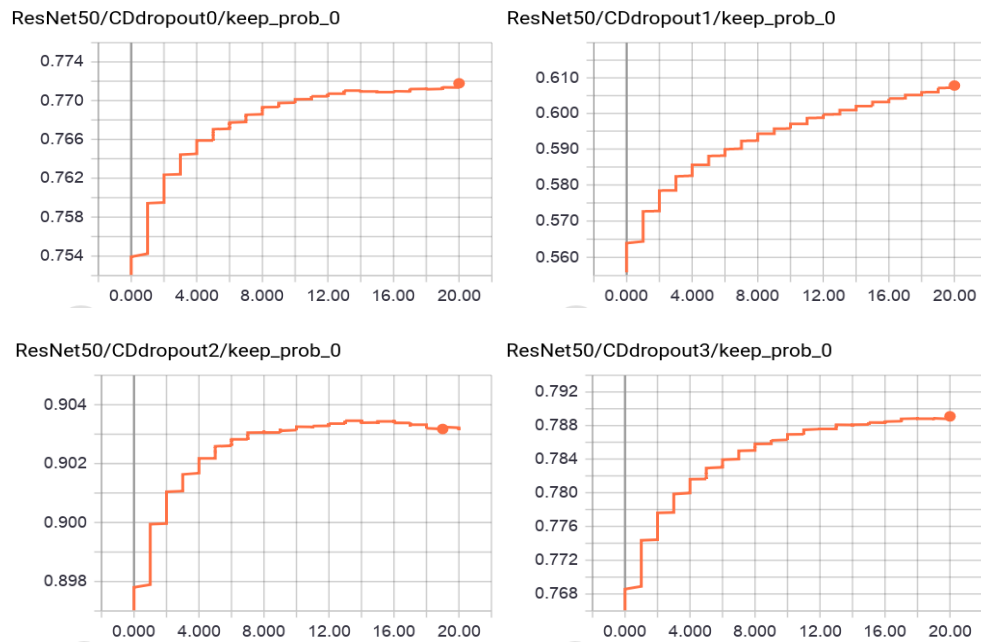


Figure 2.7: Changes along epochs of keep probability in training network with concrete dropout layer.

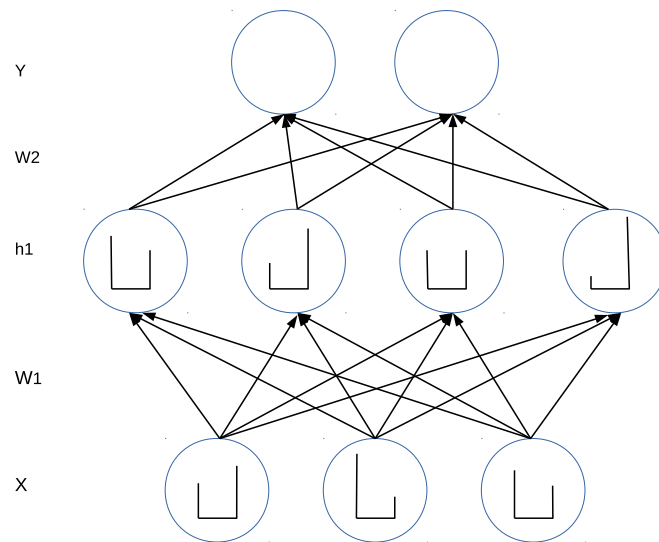


Figure 2.8: Different dropout rates for different hidden units in multi-drop.

2.2.4 Modified network architecture

After introducing dropout and concrete dropout variational inference, we will describe the modified network architecture in this subsection. The fundamental task in this work is object classification. We choose ResNet50[HZRS16] pre-trained on ImageNet as backbone for fine-tuning because of its strong ability to learn powerful representation for images. However, there is no dropout in original version of ResNet50. If we want to employ dropout variational inference to obtain reliable uncertainty estimation, dropout should be inserted into the network. In this work, we add three fully connected layer with 1024 hidden units, which are initialized from scratch, before the output layer whose dimension needs to be set to the number of classes. Then we add concrete dropout at flatten layer, and these three new added fully connected layer, respectively. There are three reasons why we modify the network in this way:

- Because we initialize major parts of network with pre-trained weights, inserting dropout in layers initialized with pre-trained weights would destroy pre-trained features. This would lead to significant drop of performance after fine-tuning.
- According to the suggestions from [SHK⁺14], insertion of dropout reduces the capacity of the model and thus a model with dropout should have larger capacity than one without dropout. Therefore we add three more fully connected layers to make sure that our model possesses large enough model capacity.
- As we have introduced in previous sub-sections, weights are major part of variational parameters. Therefore to have more weights can enhance the flexibility and capacity of approximate distribution family, which improves the quality of approximation.

From figure 2.9, we can see the sketch of our modified network architecture. More concretely, we can interpret major parts of network, which do not have dropout inserted, as a deterministic feature extractor. On the other hand, for layers with dropout inserted work as a probabilistic classifier based on aforementioned Bayesian interpretation of dropout. In training, these two parts are trained jointly in order to achieve a better balance and more optimal results. In testing, we need to marginalize possible parameters according to posterior distribution. Layers with dropout should be sampled and run multiple times to produce a predictive distribution of output class.

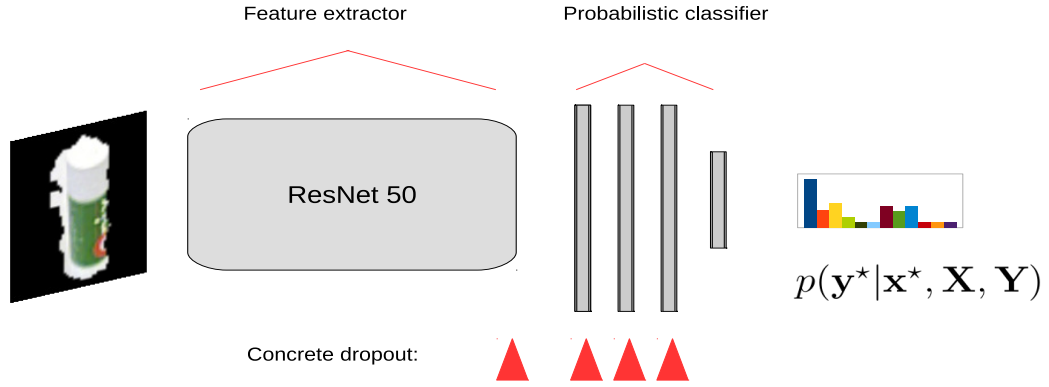


Figure 2.9: Modified network architecture.

2.3 Laplace approximation

2.3.1 Introduction

In this section, we introduce another method to approximate the real posterior distribution over weights of deep neural nets. This method is called Laplace approximation[Bis06], which is one deterministic approximation because only point estimation of model parameters instead of the learned variational parameters is required. The idea behind that is to put a Gaussian approximation on the maximum a posterior point estimate of the model parameters, which is one mode of posterior distribution. The reasons why we consider this method are as following:

- it has easy compatibility to existing network. To perform Laplace approximation, we only need one point estimation of model parameters which is already available for most of architectures. That also means, we do not need modify the training phase and we can have approximation of the true posterior for all trained models.
- it can capture relationship between model parameters. As is mentioned in the first chapter, most of approximation approaches make assumption that parameters are independent to each other for simplicity and computation burden. That could be a quite strong restrictions on the approximation leading to bad performance.

In the following, we introduce the basic idea of Laplace approximation and further introduce its scalable version for deep neural network based on Kronecker factor approximation.

Assume we have a point estimate of model parameters via maximum a posterior estimation:

$$\begin{aligned}
\boldsymbol{\omega}^* &= \operatorname{argmax}_{\boldsymbol{\omega}} \{p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X})\} \\
&= \operatorname{argmax}_{\boldsymbol{\omega}} \left\{ \frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{p(\mathbf{Y}|\mathbf{X})} \right\} \\
&= \operatorname{argmax}_{\boldsymbol{\omega}} \{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})\}
\end{aligned} \tag{2.28}$$

After taking a second order Taylor expansion of the logarithm of of posterior distribution $p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X})$ around its mode $\boldsymbol{\omega}^*$, we have:

$$\begin{aligned}
\log p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X}) &\approx \log p(\boldsymbol{\omega}^*|\mathbf{Y}, \mathbf{X}) + \\
&\quad (\boldsymbol{\omega} - \boldsymbol{\omega}^*)^T \frac{\partial \log p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X})}{\partial \boldsymbol{\omega}} + \\
&\quad \frac{1}{2} (\boldsymbol{\omega} - \boldsymbol{\omega}^*)^T \frac{\partial^2 \log p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X})}{\partial \boldsymbol{\omega}^2} (\boldsymbol{\omega} - \boldsymbol{\omega}^*) \\
&= \log p(\boldsymbol{\omega}^*|\mathbf{Y}, \mathbf{X}) - \frac{1}{2} (\boldsymbol{\omega} - \boldsymbol{\omega}^*)^T \mathbf{H} (\boldsymbol{\omega} - \boldsymbol{\omega}^*)
\end{aligned} \tag{2.29}$$

where

$$\begin{aligned}
\mathbf{H} &= - \frac{\partial^2 \log p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X})}{\partial \boldsymbol{\omega}^2} \\
&= - \frac{\partial^2 \log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}))}{\partial \boldsymbol{\omega}^2} - \frac{\partial^2 p(\boldsymbol{\omega})}{\partial \boldsymbol{\omega}^2}
\end{aligned}$$

which is negative Hessian of the log posterior. The first order term in equation 2.29 vanishes because we expand the function around a local maximum $\boldsymbol{\omega}^*$, where the first derivative is zero. If we exponentiate this equation, we can get the following Gaussian-like functional form:

$$p(\boldsymbol{\omega}|\mathbf{Y}, \mathbf{X}) \propto p(\boldsymbol{\omega}^*|\mathbf{Y}, \mathbf{X}) \exp\left\{-\frac{1}{2}(\boldsymbol{\omega} - \boldsymbol{\omega}^*)^T \mathbf{H} (\boldsymbol{\omega} - \boldsymbol{\omega}^*)\right\} \tag{2.30}$$

which means $\boldsymbol{\omega} \sim \mathcal{N}(\boldsymbol{\omega}^*, \mathbf{H}^{-1})$. Therefore, we can obtain a Gaussian approximate distribution with local maximum $\boldsymbol{\omega}^*$ as mean and inverse of negative Hessian as covariance matrix.

We assume a Gaussian prior for weights. Then it's easy to know that the second order derivative of prior distribution term $\frac{\partial^2 p(\boldsymbol{\omega})}{\partial \boldsymbol{\omega}^2}$ is a identity matrix multiplied by regularization coefficient. And the non-trivial part is the first term, second derivatives of log likelihood. To make the explanation uncluttered, we define the negative Hessian of log likelihood with $\hat{\mathbf{H}}$ instead:

$$\hat{\mathbf{H}} = - \frac{\partial^2 \log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}))}{\partial \boldsymbol{\omega}^2}$$

However, we should note that for a large training set, it's always infeasible to analyze the gradients or Hessian exactly. To resolve this, normally we estimate the expectation of gradients or Hessian in each mini-batch. That means we need to estimate Hessian with empirical average Hessian computed in mini-batches:

$$\hat{\mathbf{H}} \approx N \mathbb{E}_{p(\mathbf{Y}, \mathbf{X})}[\hat{\mathbf{H}}]$$

where N is size of training samples, and

$$\mathbb{E}_{p(\mathbf{Y}, \mathbf{X})}[\hat{\mathbf{H}}] \approx -\frac{1}{K} \sum_k \left[\frac{1}{M} \sum_i \frac{\partial^2 \log(p(y_{ik} | \mathbf{x}_{ik}, \boldsymbol{\omega}))}{\partial \boldsymbol{\omega}^2} \right] \quad (2.31)$$

where K is total number of mini-batch, and $(y_{ik}, \mathbf{x}_{ik})$ is the i -th training data sample in k -th mini-batch.

Fisher information matrix is equivalent to the expected negative Hessian of exponential family log probability:

$$\hat{\mathbf{F}} = \mathbb{E}_{p(\mathbf{X}, \mathbf{Y})}[\hat{\mathbf{H}}]$$

Therefore we can use Fisher information matrix as replacement of expected Hessian for the log likelihood term. The derivation of equivalence between expected Hessian and Fisher matrix is straightforward. We define Fisher matrix \mathbf{F} in the following:

$$\begin{aligned} \hat{\mathbf{F}} &= \mathbb{E}_{p(\mathbf{Y}, \mathbf{X})} \left[\frac{\partial}{\partial \boldsymbol{\omega}} \log p(\mathbf{Y} | \mathbf{X}, \boldsymbol{\omega}) \frac{\partial}{\partial \boldsymbol{\omega}} \log p(\mathbf{Y} | \mathbf{X}, \boldsymbol{\omega})^T \right] \\ &\approx \frac{1}{K} \sum_k \left[\frac{1}{M} \sum_i \left(\frac{\partial}{\partial \boldsymbol{\omega}} \log p(y_{ik} | \mathbf{x}_{ik}, \boldsymbol{\omega}) \frac{\partial}{\partial \boldsymbol{\omega}} \log p(y_{ik} | \mathbf{x}_{ik}, \boldsymbol{\omega})^T \right) \right] \end{aligned} \quad (2.32)$$

We include the derivation of equivalence between negative expected Hessian and Fisher matrix for case that parameter is scalar in appendix ??, which can be generalized to case for vector. Therefore we can obtain

$$\hat{\mathbf{H}} \approx N \hat{\mathbf{F}} \quad (2.33)$$

and

$$\mathbf{H} \approx N \hat{\mathbf{F}} + \tau \mathbf{I} \quad (2.34)$$

Actually, to compute outer product of gradients is feasible. But to compute the empirical average of outer product we need to save values with amount quadratic to the number of model parameters which induces large storage overhead. More than that, to inverse this

matrix, the computation complexity is cubic of dimensionality of this matrix which is the number of model parameters which is also infeasible because the number of parameters in deep neural nets can reach million easily. Therefore an efficient approximation for Fisher matrix is required.

2.3.2 Scalable Laplace approximation for neural network

In order to mitigate the computational burden in Laplace approximation above, [RBB18] proposes to use Kronecker-factored approximation curve(KFAC) in [MG15] to approximate the Fisher information matrix \mathbf{F} and perform Laplace approximation for neural network. This approximation is derived by approximating the large blocks matrix of Fisher (corresponding to each layer) by Kronecker product of two much smaller matrices. While several times more expensive to compute this, the issues of storage and inversion in Laplace approximation can be mitigated. Additionally, because we need this Fisher information matrix after training, this little expensive computation is required to perform once.

With the notations used in previous subsection, where dropout variational inference is introduced, we will derive the approximation in the following. For clarity, we repeat the definitions of existing notations and new defined ones here. A neural network transforms its input $a_0 = \mathbf{x}$ to an output $a_L = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) = \sigma(\mathbf{x}, \boldsymbol{\omega})$ through a series of L layers. The units each receive as input a weighted sum of outputs of units from the previous layer and compute their output via a non-linear activation function $\phi(\cdot)$. We denote by \mathbf{s}_i the vector of these weighted sums for i -th layer, and by \mathbf{a}_i the vector output of non-linear activation function. Computation performed for each layer $i \in 1, \dots, L$ is given as follows:

$$\mathbf{s}_i = \mathbf{a}_i \mathbf{W}_i \text{ and } \mathbf{a}_i = \phi(\mathbf{s}_i) \quad (2.35)$$

We define $\tilde{\mathbf{w}}$ to be the vector consisting of all of the network's parameters concatenated together, i.e. $\tilde{\mathbf{w}} = [\text{vec}(\mathbf{W}_1)^T \text{vec}(\mathbf{W}_2)^T \dots \text{vec}(\mathbf{W}_L)^T]^T$, where vec is the operator which vectorizes matrices by stacking their columns together. For simplicity we first define the following notations:

$$\partial v = -\frac{\partial \log p(\mathbf{y}, |\mathbf{x}\boldsymbol{\omega})}{\partial v} \text{ and } \mathbf{g}_i = \partial \mathbf{s}_i$$

The first derivatives of objective w.r.t. $\tilde{\mathbf{w}}$ is defined as follows:

$$\partial \tilde{\mathbf{w}} = [\text{vec}(\partial \mathbf{W}_1)^T \text{vec}(\partial \mathbf{W}_2)^T \dots \text{vec}(\partial \mathbf{W}_L)^T]^T \quad (2.36)$$

From equation 2.32, we know $\hat{\mathbf{F}} = \mathbb{E}[\partial \tilde{\mathbf{w}} \partial \tilde{\mathbf{w}}^T]$ which can be viewed as L by L block matrix, with the (i, j) -th block $\hat{\mathbf{F}}_{i,j}$ given by $\hat{\mathbf{F}}_{i,j} = \mathbb{E}[\text{vec}(\partial \mathbf{W}_i) \text{vec}(\partial \mathbf{W}_j)^T]$.

Noting that $\partial \mathbf{W}_i = \mathbf{a}_{i-1} \mathbf{g}_i^T$ and $\text{vec}(\mathbf{u} \mathbf{v}^T) = \mathbf{v} \otimes \mathbf{u}$, then we have

$$\text{vec}(\partial \mathbf{W}_i) = \text{vec}(\mathbf{a}_{i-1} \mathbf{g}_i^T) = \mathbf{g}_i \otimes \mathbf{a}_{i-1}$$

Then we can rewrite

$$\begin{aligned}
\hat{\mathbf{F}}_{i,j} &= \mathbb{E}[\text{vec}(\partial \mathbf{W}_i) \text{vec}(\partial \mathbf{W}_j)^T] \\
&= \mathbb{E}[(\mathbf{g}_i \otimes \mathbf{a}_{i-1})(\mathbf{g}_j \otimes \mathbf{a}_{j-1})^T] \\
&= \mathbb{E}[(\mathbf{g}_i \otimes \mathbf{a}_{i-1})(\mathbf{g}_j^T \otimes \mathbf{a}_{j-1}^T)] \\
&= \mathbb{E}[(\mathbf{g}_i \mathbf{g}_j^T) \otimes (\mathbf{a}_{i-1} \mathbf{a}_{j-1}^T)]
\end{aligned}$$

Here we need to define $\tilde{\mathbf{F}}_{i,j}$ as approximation to $\hat{\mathbf{F}}_{i,j}$ in the following:

$$\begin{aligned}
\hat{\mathbf{F}}_{i,j} &\approx \tilde{\mathbf{F}}_{i,j} \\
&= \mathbb{E}[\mathbf{g}_i \mathbf{g}_j^T] \otimes \mathbb{E}[\mathbf{a}_{i-1} \mathbf{a}_{j-1}^T] \\
&= \mathbb{E}[\mathbf{G}_{i,j}] \otimes \mathbb{E}[\mathbf{A}_{i,j}]
\end{aligned} \tag{2.37}$$

where $\mathbf{G}_{i,j} = \mathbf{g}_i \mathbf{g}_j^T$ and $\mathbf{A}_{i,j} = \mathbf{a}_{i-1} \mathbf{a}_{j-1}^T$. The expectation of a Kronecker product is, in general, not equal to the Kronecker product of expectations, and so this is indeed a major approximation to make, and one which likely won't become exact under any realistic set of assumptions, or as a limiting case in some kind of asymptotic analysis. Nevertheless, it seems to be fairly accurate in practice, and is able to successfully capture the "coarse structure" of the Fisher information Matrix. In this work, we only consider the case that layers are independent to others. Therefore the Fisher information matrix of the entire network is a diagonal block matrix, which means $\tilde{\mathbf{F}}_{i,j}$ is non-zero only for $i = j$. For i -th layer, we can compute the Fisher information matrix with equation 2.37. The two factors in Kronecker are outer product of gradients of pre-activation $\mathbf{G}_i = \mathbf{g}_i \mathbf{g}_i^T$ and outer product of activation from previous layer $\mathbf{A}_i = \mathbf{a}_{i-1} \mathbf{a}_{i-1}^T$, respectively:

$$\tilde{\mathbf{F}}_i = \mathbf{G}_i \otimes \mathbf{A}_i$$

If we treat them as covariances of resulting Gaussian. That means each Gaussian has a Kronecker factored covariance, corresponding to a matrix normal distribution[GN99], which considers the two Kronecker factors of the covariance to be the covariances of the rows and columns of a matrix. The two factors are much smaller than the full covariance and allow for significantly more efficient inversion and sampling. The final posterior of Laplace approximation for i -th layer is:

$$\mathbf{W}_i \sim \mathcal{MN}(\mathbf{W}_i^*, \mathbb{E}[\mathbf{A}_i]^{-1}, \mathbb{E}[\mathbf{G}_i]^{-1}) \tag{2.38}$$

If we consider Gaussian prior and scale of Fisher information matrix(cf.equation2.34), then the resulting posterior can be written in the following form:

$$\mathbf{W}_i \sim \mathcal{MN}(\mathbf{W}_i^*, (\sqrt{N}\mathbb{E}[\mathbf{A}_i] + \sqrt{\tau}\mathbf{I})^{-1}, (\sqrt{N}\mathbb{E}[\mathbf{G}_i] + \sqrt{\tau}\mathbf{I})^{-1}) \tag{2.39}$$

where N is the size of training set and τ is the standard deviation of Gaussian prior. However, we find that in practice, to treat them as hyper-parameters and tune them w.r.t. the predictive performance on a validation set can yield better performance. The possible

reasons for this are two fold. Firstly, nowadays the scale of dataset is really large which can exceed tens of thousand easily. The large size of data could be probably redundant and if we take the redundant data into account and thus increase N . This will lead to unreasonably high precision for Gaussian and underestimate the uncertainty significantly. Second, Laplace approximation requires positive definite Hessian and τ always serves as damp factors to fulfill this condition, therefore it needs to be chosen and tuned carefully.

Chapter 3

Conditional random field

Graphical models allow us to encode relationships between multiple variables using a concise, well-defined language. Because it's hard to solve the problems deterministically in real world. *Probabilistic graphical model*(PGM)[KFB09] can assist in this case by encoding a joint or conditional probability distribution with graph, which means that PGM can encode the relationships between multiple random variables.

The various types of PGM differ by the graph structure and the conditional independence assumptions. There are mainly two types, directed and undirected graph. Directed graphical model which is also called Bayesian network, specifies the family of distributions by means of directed acyclic graph and its joint distribution is factorized over each node which is a random variable with probability distribution conditioned on its parents. On the other hand, undirected graphical model, also known as Markov random field has undirected acyclic or cyclic graph. Its joint distribution is factorized over cliques, which is a fully-connected subgraph within the graph. This means, the marginal probability distribution in undirected graphical model is not normalized over nodes, but over cliques.

In a object recognition task, the input is set of images and output is set of predicted labels. Because the dimensionality of input is always large so that it's difficult to model it directly. Nevertheless, we want to exploit the contextual relationship between different objects. To this end, it's better to just model the relationship between output random variables, which obey the conditional distribution of output given input. *Conditional random field*(CRF), first proposed by [LMP01], is a suitable model for this purpose. In this chapter, a concise introduction of CRF is given, including definition, learning and inference in the following subsection.

3.1 Definition

CRF can model the conditional distribution of a set of output random variables \mathbf{y} given a set of inputs \mathbf{x} , that is $p(\mathbf{y}|\mathbf{x})$. In the scope of scene object recognition problem, assume we have multiple objects in one scene denoted by $\mathbf{x} = [x_1, \dots, x_n]$ and their labels $\mathbf{y} = [y_1, \dots, y_n]$. We want to learn $p(\mathbf{y}|\mathbf{x})$ with CRF and make predictions for unseen data.

To formulate the problem clearly, we denote the set of output random variables by \mathbf{y} and overall output domain by \mathcal{Y} . CRF for one scene \mathbf{x} is defined by a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where a set of nodes \mathcal{V} representing the random variables in \mathbf{y} , and a number of undirected edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ representing the dependencies between random variables. The output domain of CRF is the product of individual variable domains \mathcal{Y}_i so that we have $\mathcal{Y} = \prod_{i \in \mathcal{V}} \mathcal{Y}_i$, where $\mathcal{Y}_i = \mathcal{L}$, which is a set of all possible labels. With these notations, we can define the probability distribution of \mathbf{y} conditioned \mathbf{x} with model parameter $\boldsymbol{\theta}$ as

$$p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\mathbf{x}, \boldsymbol{\theta})} \prod_{c \in \mathcal{C}} \psi_c(y_c, x_c, \boldsymbol{\theta}) \quad (3.1)$$

where \mathcal{C} is the set of cliques of graph \mathcal{G} , and $Z(\cdot)$ is normalizer to make sure $\sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) = 1$, which is also called partition function. $\boldsymbol{\theta}$ stands for the parameters of CRF to be tuned during training phase. $\psi_c(\cdot)$ is called potential function or factor which encodes the contribution of clique c to the total density. The order of number of random variables in clique c can be arbitrary number. However, for computational reason, we only consider factor with number of random variables up to two in normal case. The first order potential is called *unary* potential which means that the clique of this kind of potential contains only one random variable, that is $\psi_i(y_i, x_i, \boldsymbol{\theta})$. Unary potential can be treated as the likelihood that a label is assigned to this random variable independently to others. The second order potential is called pairwise potential $\psi_{i,j}(y_i, y_j, x_i, x_j, \boldsymbol{\theta})$ characterized by edges in the graph. This kind of potential states the compatibility of two random variables (y_i, y_j) being assigned to a certain pair of categories.

To make the explanation more clear, we illustrate a simple example in figure 3.1, we ignore $\boldsymbol{\theta}$ here and explain it later and use $\psi_{i,j}$ for $\psi_{i,j}(y_i, y_j, x_i, x_j)$ to make expression of conditional probability uncluttered:

$$\begin{aligned} p(y_1, y_2, y_3, y_4) &= \frac{1}{Z(x_1, x_2, x_3, x_4)} \psi_1 \psi_2 \psi_3 \psi_4 \psi_{1,4} \psi_{1,2} \psi_{2,3} \psi_{2,3} \psi_{1,2,4} \\ &= \frac{1}{Z(x_1, x_2, x_3, x_4)} \psi_1 \psi_2 \psi_3 \psi_4 \psi_{1,4} \psi_{1,2} \psi_{2,3} \psi_{2,3} \end{aligned}$$

where the transformation from first line to second line comes from the ignorance of potentials with order higher than two.

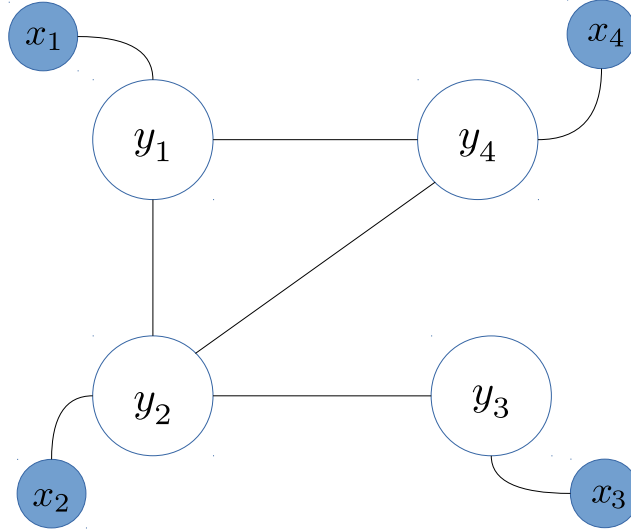


Figure 3.1: Simple example of conditional random field.

According to Hammersley-Clifford theorem [HC68], the factors must be always positive, we can employ log-linear models for the potential functions, which results in:

$$\begin{aligned}
 p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) &= \frac{1}{Z(\mathbf{x}, \boldsymbol{\theta})} \prod_{c \in \mathcal{C}} \exp(\langle \phi(x_c, y_c), \boldsymbol{\theta} \rangle) \\
 &= \frac{1}{Z(\mathbf{x}, \boldsymbol{\theta})} \exp\left(\sum_{i \in V} \langle \phi_u(x_i, y_i), \boldsymbol{\theta}_u \rangle + \sum_{(i,j) \in \mathcal{E}} \langle \phi_p(x_i, x_j, y_i, y_j), \boldsymbol{\theta}_p \rangle\right)
 \end{aligned} \tag{3.2}$$

being $\langle \cdot, \cdot \rangle$ the inner product, and $\phi(x_c, y_c)$ the sufficient statistic of the factor over clique c . Because we consider the potentials up to second order, only feature function for unary potential $\phi_u(\cdot)$ and feature function for pairwise potential $\phi_p(\cdot)$ are considered here. This function can also be treated as feature function which extracts feature of the data over clique.

Now we have made assumptions of graph structure and conditional dependency for the graph. This can be interpreted as a filter of a probability distribution family. Only the class of probability distribution that satisfies these assumptions such as up to second order potentials can be expressed and learned with this graph. To further identify different probability distribution in this class, we need to parameterize it and identify it with specific parameters which can be learned during training.

Parametrization of the probability distribution is expressed explicitly in Eq.3.2. This means that the dimension of the parameter for each kind of potential should be the same as

the dimension of feature function. In this work, the feature function should be a function of input and class label, while the feature function can be the function of only input and weight should be defined as matrix whose output dimension is the number of class serving as a linear mapping. In this work, we only consider the former case and utilize the probability vector $p(y_i|x_i)$ which is output from Bayesian neural network as unary potential and binary co-occurrence matrix $\mathbf{H}(y_i, y_j)$ as pairwise potential. The details of these two potentials will be explained in experimental part. Based on these two kinds of potentials, we have the formulation of CRF model for one scene as follows:

$$p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\mathbf{x}, \boldsymbol{\theta})} \exp(\theta_u \sum_{i \in V} p(y_i|x_i) + \theta_p \sum_{(i,j) \in \mathcal{E}} \mathbf{H}(y_i, y_j)) \quad (3.3)$$

where $\boldsymbol{\theta} = \{\theta_u, \theta_p\} \in \mathbb{R}^2$.

3.2 Learning

The most popular approach is estimate $\boldsymbol{\theta}$ is maximum likelihood estimate(MLE). The likelihood for one scene is defined in Eq.3.3. Normally we have large size of scenes in training set where each of them contains multiple objects. We denote training set by $\mathcal{D} = [d_1, \dots, d_m]$, where each scene is represented by $d_i = (\mathbf{x}^i, \mathbf{y}^i)$. For each scene, we create a CRF model for it. Therefore we can denote the whole set of graph by $\mathcal{G} = [\mathcal{G}_1, \dots, \mathcal{G}_m]$, where $\mathcal{G}_i = (V_i, \mathcal{E}_i)$. This is consistent to previous notations when we make little modifications for \mathbf{x} and \mathbf{y} : $\mathbf{x}^i = [x_1^i, \dots, x_n^i]$ and $\mathbf{y}^i = [y_1^i, \dots, y_n^i]$. For computational convenience and stability, negative log likelihood as objective function is always employed. The objective in MLE on training set is defined as follows:

$$\mathcal{NLL}(\boldsymbol{\theta}; \mathcal{D}) = - \sum_{i=1}^m (\theta_u \sum_{j \in V_i} p(y_j^i|x_j^i) + \theta_p \sum_{(j,k) \in \mathcal{E}_i} \mathbf{H}(y_j^i, y_k^i) - Z_l(\mathbf{x}^i; \boldsymbol{\theta})) \quad (3.4)$$

where

$$Z_l = \log \left(\sum_{\mathbf{y}^i \in \mathcal{Y}} (\theta_u \sum_{j \in V_i} p(y_j^i|x_j^i) + \theta_p \sum_{(j,k) \in \mathcal{E}_i} \mathbf{H}(y_j^i, y_k^i)) \right) \quad (3.5)$$

With this expression, our problem is defined in the following form with MLE:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} [\mathcal{NLL}(\boldsymbol{\theta}; \mathcal{D})] \quad (3.6)$$

The log likelihood function of CRF is proved to be concave [KFB09], which means the local optimizer is also the global optimizer although it's not unique. To optimize this objective w.r.t. the model parameters on dataset with large size, we can employ stochastic gradient descent (SGD), which requires to calculate the gradient of object w.r.t. the model

parameters $\boldsymbol{\theta}$. The gradients can be computed as follows, where we use $\boldsymbol{\theta}$ to represent parameters in order to make the expression clear:

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{NLL}(\boldsymbol{\theta}; \mathcal{D}) = - \sum_{i=1}^m (\phi(\mathbf{x}^i, \mathbf{y}^i) - \mathbb{E}_{p(\mathbf{y}|\mathbf{x}^i; \boldsymbol{\theta})} [\phi(\mathbf{x}^i, \mathbf{y})]) \quad (3.7)$$

where the expectation $\mathbb{E}_{p(\mathbf{y}|\mathbf{x}^i; \boldsymbol{\theta})} [\phi(\mathbf{x}^i, \mathbf{y})]$ stands for the partial derivative of the log-partition function, which is retrieved by:

$$\mathbb{E}_{p(\mathbf{y}|\mathbf{x}^i; \boldsymbol{\theta})} [\phi(\mathbf{x}^i, \mathbf{y})] = \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}|\mathbf{x}^i; \boldsymbol{\theta}) \phi(\mathbf{y}, \mathbf{x}^i) \quad (3.8)$$

The gradients of parameters corresponds to the difference between the empirical expectation of its associated sufficient statistics (the sufficient statistics of ground truth) and its expectation over the estimated probability distribution $p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$. The goal of the optimization is to decrease the gradients until they reach zero, which means that we match of first moment of sufficient statistics over the output distribution. Therefore this way is also called moment matching.

When taking a closer look at Eq. 3.7 and Eq. 3.8, to calculate the gradients in each step of SGD, we need to compute $p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$, which requires to compute the partition function $Z(\cdot)$. Because the number of possible assignments grows exponentially in the number of classes, it's infeasible to get a exact solution in practice. This problem also appears when we want to obtain a probability distribution (marginal inference) instead of a single prediction (MAP inference) in testing time. To compute the probability distribution in probabilistic modeling always refers to as inference. There are a bunch of research works focusing on inference. Because the CRF model in this work is fully connected, we adopt a popular and effective approach which is called loopy belief propagation and introduced in next subsection. This approach can provide a approximation to the log-partition function $Z_l(\cdot)$, which is required in 3.8 and to the marginal probability distribution of random variables in \mathbf{y} to be plugged in 3.7 in the graph with cycles and exact solution in graph without cycle.

3.3 Inference

Loopy belief propagation introduced by [Pea14], also known as sum-product algorithm, is a popular approach for probabilistic inference in graphical models. Briefly, this approach is based on the exchange of statistical information among the nodes in the graph according to their relationships. The core idea behind that is for each node in the graph to maintain its belief, which is a k -dimensional vector $\mathbf{b}(y_i)$ if the node is discrete random variable where k is the number of class. The i -th entry of this vector indicates the probability that i -th class is assigned to this node. The belief of a node evolves as it receives *messages* from it

neighbors. A message $\mathbf{m}_{ij}(y_j)$ sent from node y_i to node y_j encodes belief of y_i about what class node y_j should belong to. This message, in turn, based on the messages it received from all the other neighbors except for the recipient in a previous time-step.

The algorithm keeps sending messages until the graph is calibrated, which means that the messages of two consecutive iterations are less than one threshold. However, in practice, the algorithm can not calibrate or converge. A maximum number of iteration can be set to obtain satisfactory results. The messages computation is as follows:

$$\mathbf{m}_{ij}(y_j) = \sum_{y_i \in \mathcal{V}_i} \psi_{ij}(y_i, y_j, x_i, x_j, \theta_p) \psi_i(y_i, x_i, \theta_u) \prod_{y_k \in \mathcal{N}(y_i) \setminus y_j} \mathbf{m}_{ki}(y_i) \quad (3.9)$$

where $\mathcal{N}(y_i) \setminus y_j$ is the set of node neighbors of y_i except for y_j . Once the iteration reaches maximum number or the graph is calibrated, we can compute the belief of each node as follows:

$$\mathbf{b}(y_i) = \frac{1}{Z_i} \psi_i(y_i, x_i, \theta_u) \prod_{y_j \in \mathcal{N}(y_i)} \mathbf{m}_{ji}(y_i) \quad (3.10)$$

where Z_i is the normalizer of unnormalized belief of node y_i to ensure the sum of marginal probability of y_i is 1.

For the implementation, we employ the library in C++ from [RSGGJ15]. However, the weights in this library are assumed to be category-independent, which means that the weights are defined as weight matrix to build a linear mapping from feature function to potential function. Therefore we need to modify the type of weights in this library in order to incorporate different type of weight defined in Eq. 3.3, which serve as coefficient of potential function instead of a linear mapping inside the potential function. To note that the messages are initialized uniformly during loopy belief propagation.

Chapter 4

Experiments and Discussion

4.1 Introduction

In this chapter, we perform different experiments for the following purposes:

- to evaluate the performance of uncertainty estimation of dropout variational inference and its variants as well as Laplace approximation. Comparisons and analysis towards the results of these two kinds of approach are given.
- to evaluate performance of training a domain specific classifier, with different strategies for collecting training data including manual labeling and automatic labeling, where the latter one is chosen based on uncertainty estimation.
- to evaluate performance of classifier including context information via CRF.

Before looking into the results, we need to specify the details of experiments such as the dataset, evaluation metrics. Besides, the detailed parameters of the model will be reported in section of each experiment to avoid confusions. Models were implemented in Tensorflow[ABC⁺16], and optimization was performed using RMSprop with initial learning rate of $1e^{-5}$ and L2 regularization with coefficient of $3.5e^{-6}$ as well dropout regularization with coefficient of $1.0e^{-5}$. We used early stopping with all methods, where the amount of epochs to run was determined based on performance on a validation set. We set the number of maximum epoch as 20.

4.1.1 Dataset

WRGBD[LBRF11] This dataset is a large-scale dataset of 300 household objects captured from multi-viewpoint. The objects are organized into 51 categories, some of categories and their subtrees. To note that category level recognition involves classifying objects with similar semantic appearance such as cereal boxes with different texture into

same category, instead of just physical appearance. Instance level recognition is to classify objects with similar physical appearance to the same category. Therefore the overlapping of features between classes in category recognition is larger than that in instance recognition. More abstract and semantic concepts are expected to be learned in category recognition. Each object was placed on a turn table and captured from a systematically sampled view hemisphere with 15° step in elevation (from 30° to 60°) and 2° step in azimuth (from 0° to 360°) (cf. figure 4.1). The total size of entire dataset is around 160.9×10^3 .

UniHB To simulate application-specific situation in deployment of robots, we recorded a similar dataset by putting objects on a turn table in front of the robot and recording partial views of objects. We follow the same methodology as [LBRF11] suggests for WRGBD, but with only one object instance for each category in the dataset. This analogue of WRGBD dataset is called the IAI-ODU dataset of UniHB. The size of data with elevation 45° is around 8.6×10^3 and that of data with elevation 30° and 60° is around 17.1×10^3 .

While the UniHB dataset’s setup strives to mimic the WRGBD one, the differing capturing conditions such as different equipment and light conditions even appearances of objects, result in obvious changes in feature space, thus yielding a significant drop on classification accuracy. In addition to the existing 51 categories, there are 28 novel objects(cf. figure 4.2) that do not belong to the 51 categories and these objects are treated as out-of-distribution(OOD) data for testing uncertainty estimation of the model. To note that data of these novel objects are recorded in a slight different way, that is, their view points are sampled with 15° step in elevation (from 30° to 60°) and 5° step in azimuth (from 0° to 360°). The size of this OOD dataset is around 6.0×10^3 .

T-LESS[HHO⁺17] This dataset features thirty commodity electrical parts which have no significant texture, discriminative color or distinctive reflectance properties, and often bear similarities in shape and/or size. Furthermore, a unique characteristic of the objects is that some of them are parts of others. All images were obtained with an automatic procedure that systematically sampled images from a view sphere, resulting in around 38.0×10^3 training and 10.0×10^3 test images. The training images depict objects in isolation with a black background, while the test images originate from twenty table-top scenes with arbitrarily arranged objects. The training images of each object are captured from a systematically sampled view sphere with 10° step in elevation (from 85° to -85°) and 5° step in azimuth. In this work we only consider the cropped images of objects in testing scenes, resulting a testing set with size 69.5×10^3 , while the scene context information can be exploited with conditional random field in the final experiment.

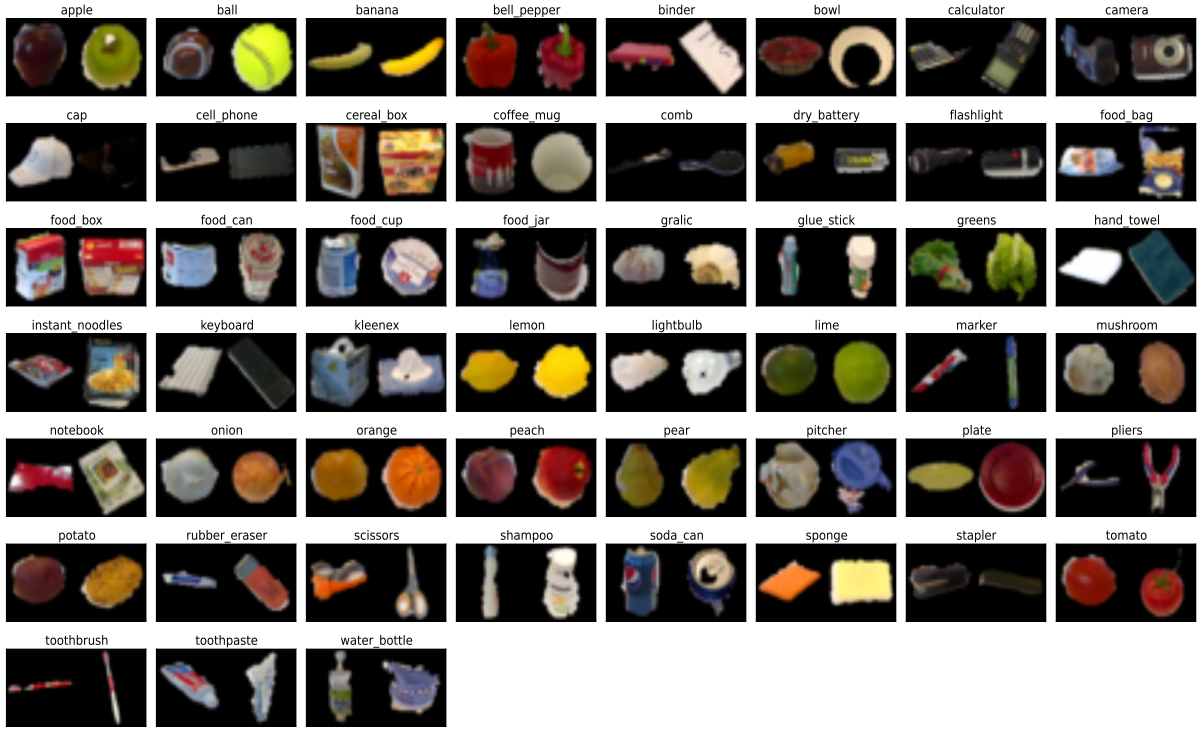


Figure 4.1: Example of masked images of objects from 51 categories in WRGBD and UniHB dataset. In each category, the left is from WRGBD and the right is from UniHB. We randomly pick one instance for the objects of WRGBD. We can see some light and appearance difference between objects in these two datasets.

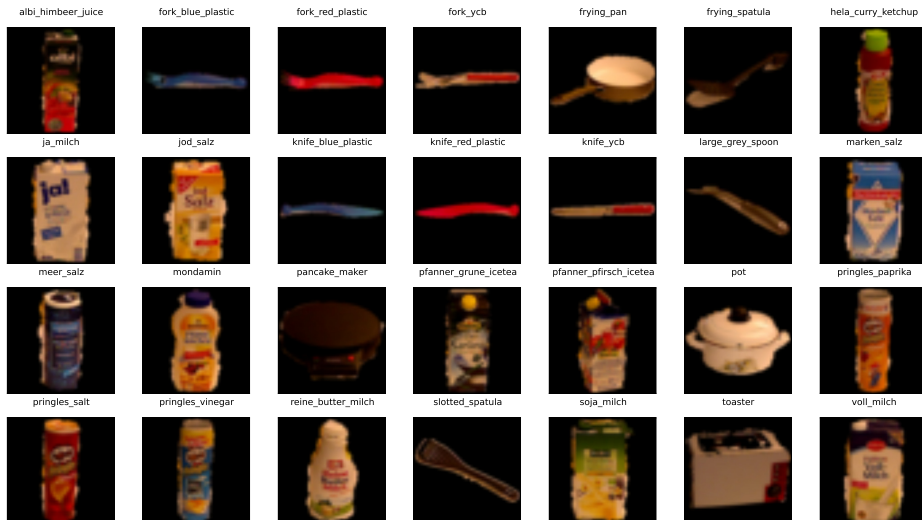


Figure 4.2: Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples.

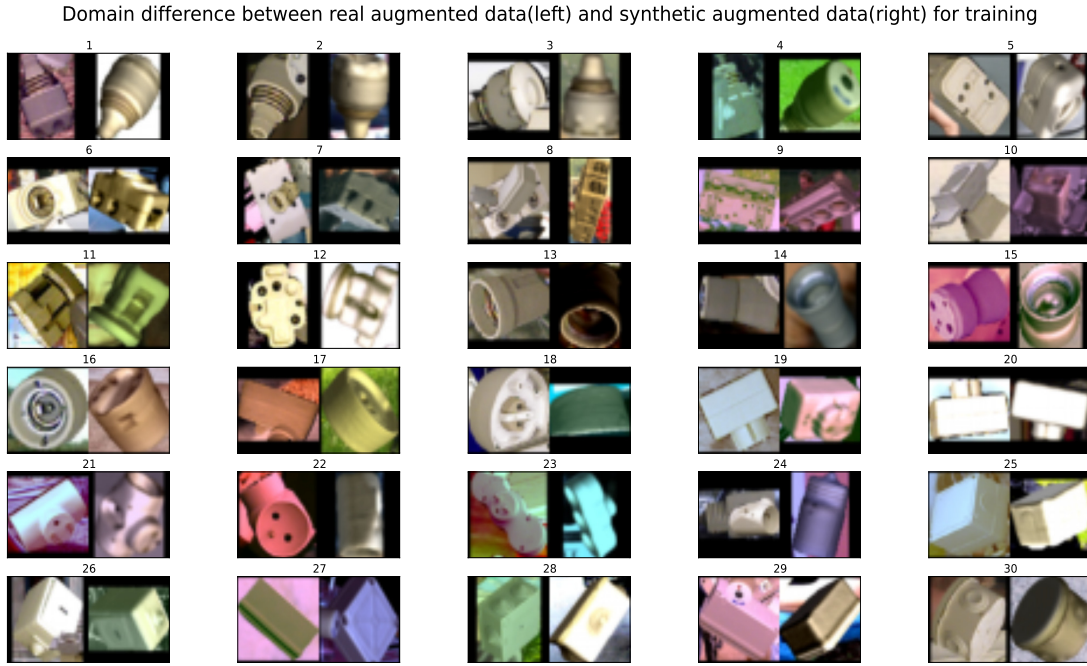


Figure 4.3: Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples.

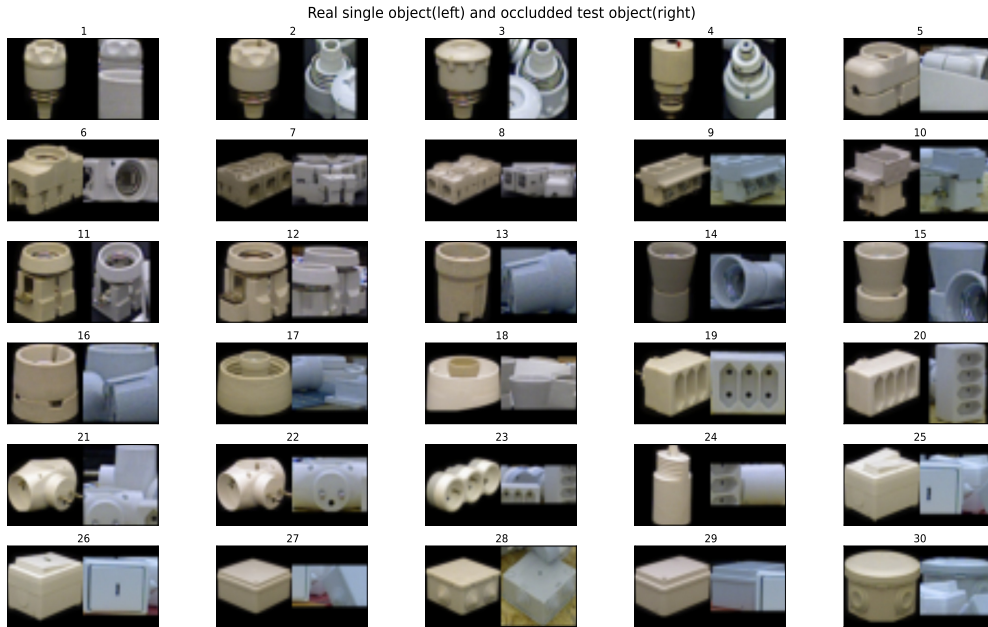


Figure 4.4: Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples.

4.1.2 Uncertainty measures

For each prediction, we can obtain one predictive probability distribution from our model with equation 2.23. In order to quantify uncertainty of the prediction, there are different metrics to measure the uncertainty of prediction, which are introduced in the following. We define x^* as one test data sample and \mathcal{D} as our training set.

Confidence is defined as the maximum probability of the output predictive distribution, whose index is the class prediction.

$$conf = \max_c [p(y = c|x^*, \mathcal{D})] \quad (4.1)$$

where $conf \in [0, 1]$, $c \in \mathcal{P} = \{0, \dots, |\mathcal{P}| - 1\}$, and \mathcal{P} represents output space in which label is expressed as number index which is transformed into one-hot representation in computation of objective function. The larger this quantity is, the less uncertain the prediction is.

Predictive entropy is quantity that captures the average amount information contained in predictive distribution[Sha48]:

$$\mathcal{H}[p(y|x^*, \mathcal{D})] = - \sum_c p(y = c|x^*, \mathcal{D}) \log p(y = c|x^*, \mathcal{D}) \quad (4.2)$$

where c is the possible class y can take, $\mathcal{H}(\cdot) \in [0, \log|\mathcal{P}|]$, the larger this quantity is, the more uncertain the prediction is.

Mutual information Mutual information between the prediction y and the weights posterior offers a different uncertainty measure when compared with the aforementioned ones. This measure is widely used in active learning tasks[HHGL11], which is also called Bayesian active learning disagreement (BALD). The definition of this measure is as follows:

$$\begin{aligned} \mathcal{I}[y, \omega|x^*, \mathcal{D}] &= \mathcal{H}[y|x^*, \mathcal{D}] - \mathbb{E}_{p(\omega|\mathcal{D})} [\mathcal{H}[y|x^*, \omega]] \\ &= - \sum_c p(y = c|x^*, \mathcal{D}) \log p(y = c|x^*, \mathcal{D}) \\ &\quad + \mathbb{E}_{p(\omega|\mathcal{D})} [- \sum_c p(y = c|x^*, \omega) \log p(y = c|x^*, \omega)] \\ &\approx - \sum_c p(y = c|x^*, \mathcal{D}) \log p(y = c|x^*, \mathcal{D}) \\ &\quad + \mathbb{E}_{q(\omega)} [- \sum_c p(y = c|x^*, \omega) \log p(y = c|x^*, \omega)] \end{aligned} \quad (4.3)$$

where $\mathcal{I}(\cdot) \in [0, \log|\mathcal{P}|]$. This quantity considers the effect of approximate posterior distribution more directly. Compared with aforementioned uncertainty measures, this one should be able to capture the model uncertainty more accurately. To think about it intuitively, this quantity measures the information gain between the entropy of predictive output distribution and the expected entropy of output distribution wr.t. weights posterior. It will be low only if the predictive distribution agrees with most of possible models(weights realizations), which means that the model is sure about its prediction. Otherwise, it will be high because most of possible models do not agree with the other models and thus the predictive distribution will be more uniform and thus has higher entropy.

4.1.3 Evaluation metrics

As is stated in [GBR07], the goal of probabilistic prediction is to maximize the sharpness of predictive distribution subject to calibration. Calibration refers to the statistical consistency between the predictive probability and the occurrence of observations, which is the frequency of the event. Therefore we employ different metrics including both the accuracy and other quantities related to calibration as well as summary of accuracy and calibration. Additionally, we also employ histogram and diagram to express the results visually. While the visual one can show us the results more intuitively, the quantitative one can allow us to evaluate the results more objectively. The comparison between them may help us to examine if visual metrics are corresponding to the numerical metrics, which may provide more insights in evaluating uncertainty estimation.

Uncertainty histogram is a intuitive visual tool for analyzing the statistics of the uncertainty estimation. Compared with normal histogram, there is one difference to stress on. In order to make visual effect more clear and hence the analysis easier, **normalizer** of each type of prediction is the size of this type of predictions instead of size of all predictions. In detail, we have three types of prediction for plotting the histogram, which are:

- **correct prediction**
- **miss-classification**
- **out-of-distribution(OOD)**

The range of y axis is $[0, 1]$ because we normalize the number in each bin, and range of x axis is set by the range of corresponding type of uncertainty measure.

Reliability diagram(Calibration curve) is another visual tool for expressing **calibration** performance of the model[GPSW17], which plots the frequency of success(accuracy of predictions in specific bin) as a function of confidence, which is predictive likelihood of prediction. If the model is perfectly calibrated, this function should be overlapping with

the diagonal line exactly. The closer this curve is to diagonal curve, the better the calibration performance is. In order to draw the curve, we firstly group the predictions into M interval bins w.r.t. confidence and then calculate the accuracy of predictions in each bin. We use $M = 20$ in this work. To quantify the proximity to the diagonal curve, we use two metrics:

- **Expectation calibration error (ECE)**: in order to obtain a more objective measure of calibration quality, we can compute the **expected calibration error** by computing the weighted average of difference between accuracy and confidence:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} |acc(B_m) - conf(B_m)| \quad (4.4)$$

where n is the number of samples, and $acc(B_m)$ represents the accuracy of samples, $conf(B_m)$ the average predicted confidence of samples in m -th bin. We can see that this metric measures the inconsistency between the statistics and predictive distribution. To note that this metric only considers the calibration quality instead of accuracy. One drawback of this metric is that, because the weights are computed based on number of samples in each bin, the metric would bias if most of predictions are clustering in few bins. Then the weights of these bins are much larger than other bins and other bins with large error would be ignored.

- **Maximal calibration error (MCE)**: Considering the drawback mentioned in ECE, we make use of this metric. Additionally, in high risk applications where reliable confidence measures are absolutely necessary, we may wish to take the worst case into account. It is defined as follows:

$$MCE = \max_m |acc(B_m) - conf(B_m)| \quad (4.5)$$

Proper scoring rules: Scoring rules provides a **summary** measure in the evaluation of probabilistic forecasts. It assigns a numerical score based on the difference between predictive distribution and real distribution.

Because we want to make predictions for the future and also have a suitable measure of the uncertainty associated with them (see [GR07] for a review). We can define the scoring rule as a function $\mathcal{S}(p(y|\mathbf{x}), (y|\mathbf{x}))$ that evaluates the quality of predictive distribution $p(y|\mathbf{x})$ relative to an event $y|\mathbf{x} \sim q(y|\mathbf{x})$ where $q(y|\mathbf{x})$ represents the true distribution over $(y|\mathbf{x})$. Consequently, the expected scoring rule is:

$$\mathcal{S}(p, q) = \int q(y|\mathbf{x}) \mathcal{S}(p, (y|\mathbf{x})) dy \quad (4.6)$$

$\mathcal{S}(p, q)$ is proper if $\mathcal{S}(p, q) \leq \mathcal{S}(q, q)$, with equality holds if and only if $p(y|\mathbf{x}) = q(y|\mathbf{x})$. Here we adopt two simple and famous proper scoring rules in which the less it is, the better the performance is:

- **averaged negative log likelihood(NLL)**:NLL is a popular metric for evaluating predictive uncertainty [QCRS⁺05]. This metric considers aforementioned confidence as likelihood. The smaller this metric is, the better the predictive distribution is.

$$NLL = -\frac{1}{|\mathcal{D}_{test}|} \sum_{i=1}^{|\mathcal{D}_{test}|} \log(p(y_i = c_i|\mathbf{x}_i)) \quad (4.7)$$

where c_i is the ground truth label for \mathbf{x}_i .

- **Brier score** is the mean squared error between the target distribution(one-hot encoding label) and predictive distribution:

$$BS = -\frac{1}{|\mathcal{D}_{test}|} \sum_{i=1}^{|\mathcal{D}_{test}|} (\mathbf{y}_i^{gt} - p(\mathbf{y}_i|\mathbf{x}_i))^2 \quad (4.8)$$

where \mathbf{y}_i^{gt} is the one-hot encoding ground truth label for \mathbf{x}_i .

Separability metrics: In addition to the summary metrics of predictive probability. We are also interested in the separability between different types of predictions, which can assist the down-stream tasks if they are highly separable (e.g. separate correct predictions and false predictions or out-of-distribution data to improve robustness of system and provide more safety guarantee). To note that, we treat the correct predictions as positive samples in the two following metrics.

- **Area under Receiver Operating Characteristic curve(AUROC)**:Since one of our goals is to choose automatically labeled data based on uncertainty estimation, it's necessary to evaluate the separability between correct predictions and false predictions or out-of distribution data predictions. ROC curve describes the relationship between true positive rate($tpr = \frac{tp}{tp+fn}$) and false positive rate($fpr = \frac{fp}{fp+tn}$). Moreover, AUROC can be interpreted as the probability that a positive samples has a greater score than a negative samples. However, the drawback of AUROC is that, the normalizers of two kinds of rate in ROC curve are independent to each other. When these two normalizers differs too much. AUROC can provide misleading conclusion. For example, if the number of negative predictions is much higher than positive one, the AUROC could still achieve a relatively high value, while in fact there are already many false positives.
- **Area under Precision Recall curve(AUPR)**: Considering the downside of AUROC, another evaluation metric is employed, that is AUPR. AUPR describes the

relationship between precision($pr = \frac{tp}{tp+fp}$) and recall($tpr = \frac{tp}{tp+fn}$), which resolves the problem of different base number. In the previously mentioned case, though the AUROC is high, the AUPR will be low because the precision is low.

4.2 Uncertainty estimation experiments

4.2.1 Experiments I

In this experiments, we want to evaluate our model on a relatively easy task which is instance recognition at first. That means, the appearances of different classes as well as the out-of-distribution data are highly discriminable. To this end, we separate WRGBD dataset into two subsets based on the instance class. Subset I contains objects with instance class from 0 to 199 (assuming that here we use index to denote the instance label) and Subset II contains objects with instance class from 200 to 299.

We train our model with objects captured in elevation 30° and 60° of Subset I (in which we split off 20% of training set as validation set for model selection in training). The size of training set and validation set is around 71.0×10^3 . Then we test the model on objects captured in elevation 45° of both Subset I and Subset II. In this experiment, the objects in Subset II serve as out-of-distribution samples because they are not present during training. The size of test set is around 35.6×10^3 and that of OOD dataset is around 17.8×10^3 .

In this experiment, we just want to test and compare the performance of different approaches. Because the OOD data is highly different to the training set. The model should express high uncertainty when facing these data. Nevertheless, the miss-classifications should have higher uncertainty than the correct predictions because the model should be confident when it makes predictions correctly. Otherwise the model is overconfident when it's confident with its incorrect predictions, which can induce hazardous consequence in some safety-critical applications. As we can see in the figure 4.5, the original ResNet50 model is highly overconfident, which assigns high confidence to nearly 50% of miss-classification and nearly 30% of OOD data in the highest bin. While the model with concrete dropout and multi-drop can significantly decrease these two proportions. As we can see, different uncertainty show similar trend in this experiment.

In addition to uncertainty histogram, we also plot the calibration curve of each model evaluated on both test set and OOD dataset in figure 4.6. As expected, the calibration performance of original model is much worse than those of the other two models. This result is corresponding to the results in [GPSW17], that the model is highly overconfident. And the concrete dropout and multi-dropout can mitigate the undesired problem significantly.

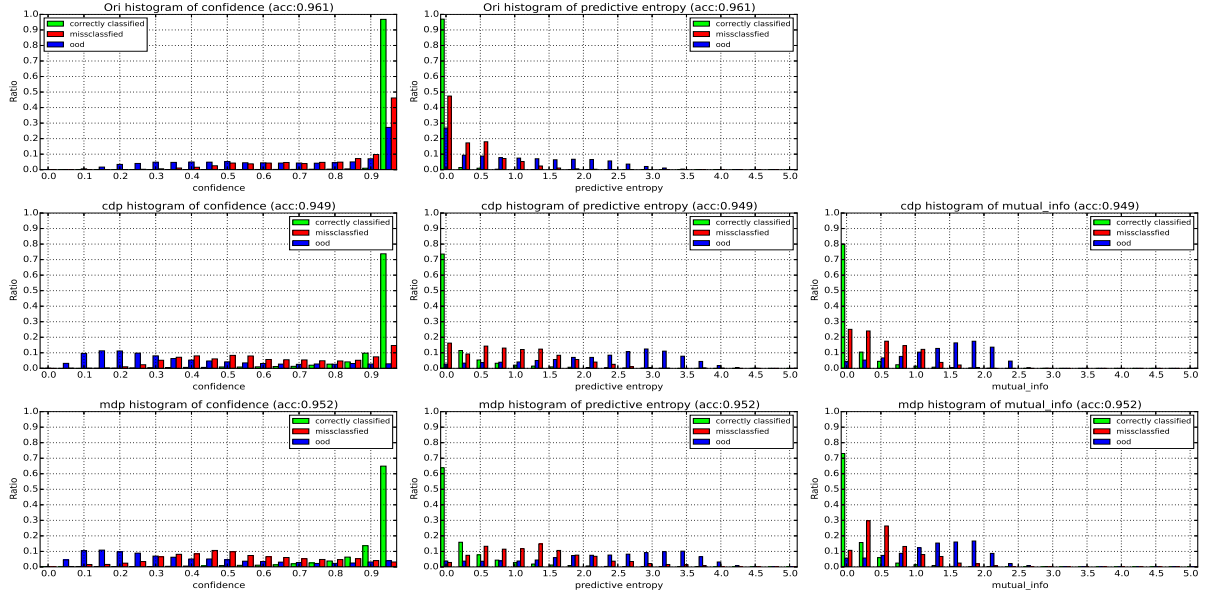


Figure 4.5: Uncertainty(confidence, predictive entropy, mutual information) histograms of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout.

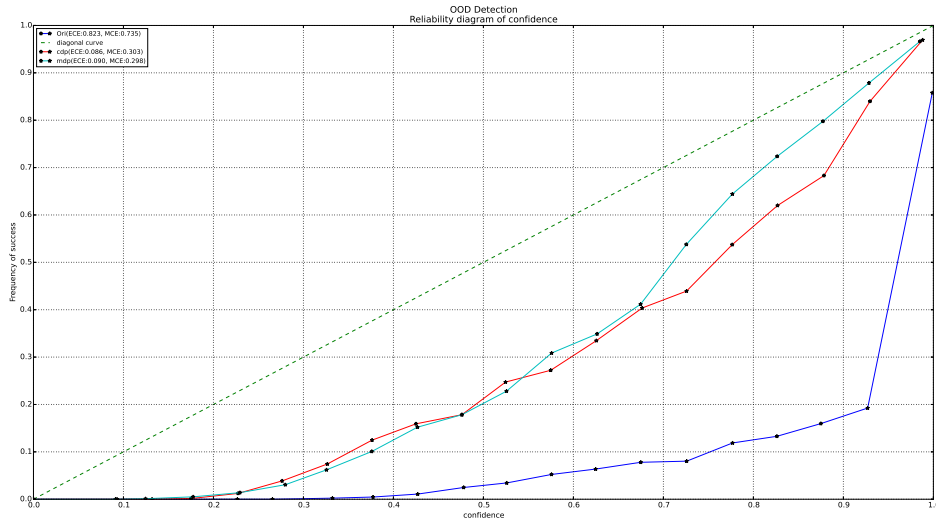


Figure 4.6: Calibration curves of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout.

4.2.2 Experiments II

In this experiment, we evaluate our model on a category recognition task. The difficulties are expressed in two-fold. Firstly, in this task, the model is confronted with classes with more overlapping. Second, because we want to simulate situation of deploying robots in real world scenario and facing data in real world, the UniHB dataset with domain gap(cf. figure4.1) to WRGBD is used to achieve this goal. This means, the uncertainty estimation should not only be able to perform well on dataset with exactly same distribution to training set, but also to generalize well to dataset with domain gap which may decrease the accuracy significantly.

Accordingly, we use the entire WRGBD dataset including all view points, whose size is around 160.9×10^3 , to train our model (in which we split off 20% of training set as validation set for model selection during training). When it comes to UniHB dataset, we treat objects captured in elevation 30° and 60° in this dataset as **adaptation set**, whose size is around 17.1×10^3 . We test performance of uncertainty estimation on this adaptation dataset. The images of 45° , whose size is around 8.6×10^3 , are used for final evaluation after we fine-tune the model with subset of adaptation set to obtain a domain specific model. The latter part will be experimented in next section. Besides, we also evaluate the uncertainty estimation on OOD dataset(cf. figure 4.2) whose size is around 6.0×10^3 .

We firstly explain protocols in the following for different kinds of metric, which can help understanding the plots and extracting useful information more easily and quickly.

- The following visual metrics are chosen in one of three runs of different random seeds. The average quantitative results are given in tables following the visual metrics.
- The calibration curve with title "Classification" on the left is plotted **only** on predictions of test set. The one with title "OOD detection" on the right is plotted on predictions of **both** OOD dataset and test dataset.
- The ROC curve and PR curve measure the separability between two types of prediction. The ones with title "classification" on the left measure the separability between **correct prediction** and **miss-classification**. The ones with title "OOD detection" on the right measure the separability between **correct prediction** and **OOD prediction**. In all ROC curve and PR curve, correct prediction is always chosen to be positive.
- As is shown in uncertainty histogram, we have used three uncertainty measures. Each one has its own ROC curve and PR curve of correct prediction versus miss-classification or OOD prediction. In the following plots of ROC curve and PR curve, we only show the one with highest area under curve. In detail, **Confidence** is chosen in case of correct prediction versus miss-classification. **Mutual information** is chosen in case of correct prediction versus OOD prediction.

Comparison with Ensemble

In this subsection, we will show the results in not only a qualitative (visual) way, but also in a quantitative way. The approaches we compare here include original version of ResNet50(ori), modified ResNet50 with concrete dropout(cdp), modified ResNet50 with multiple dropout(mdp) as well as their ensemble version. The members in ensemble are just initialized from different random seeds and no other techniques for enhancing the performance are used. In order to quantify the results more objectively, we average the results from three different random seeds and report the mean and the standard deviation.

Uncertainty histogram In the following figure 4.7, we can see uncertainty histograms for different approaches and different uncertainty measures. Firstly, let's compare them vertically. We can see the improvement visually. The original version is still overconfident, which express low uncertainty to more than 50% of OOD data and nearly 40% of miss-classification. The cdp and mdp can lower this percentage to around 10%, although at the same time the percentage of correct predictions is decreased. But we can know that the predictions for which the model expresses low uncertainty are more likely to be correct. And the ensemble of cdp and mdp can perform even better with miss-classification. However, the ensemble of mdp does not work as well as cdp version on OOD data, which can be seen in the histogram on the last row of the figure. The reason for this will be investigated in a ablation study later. Then we can compare them horizontally, it's shown that different uncertainty measures have similar trends in the histogram. Therefore we need other quantitative metrics for them in order to evaluate different approaches as well as metrics more accurately.

Reliability diagram In the figure 4.8, on the left it's figure of calibration curves of different approaches, plotted from predictions on test dataset. We can see that the calibration performance is improved a lot by the cdp and mdp. Additionally, their ensemble versions can give even better results which are tightly overlapping with the diagonal curve. On the right there are curves plotted from predictions on both test dataset and OOD dataset. Similar ranks and trends as the left curves appear here. After adding the OOD data into evaluation, all of the calibration curves are pulled down, which means that the calibration performance get worse. If all OOD data is assigned with low confidence, then curves between left and right should be similar because the accuracy in middle or high confidence interval would not change a lot or be pulled down a lot. By comparison these two figures, we can also see that the robustness against OOD data of cdp, mdp and their ensembles are better than the original version.

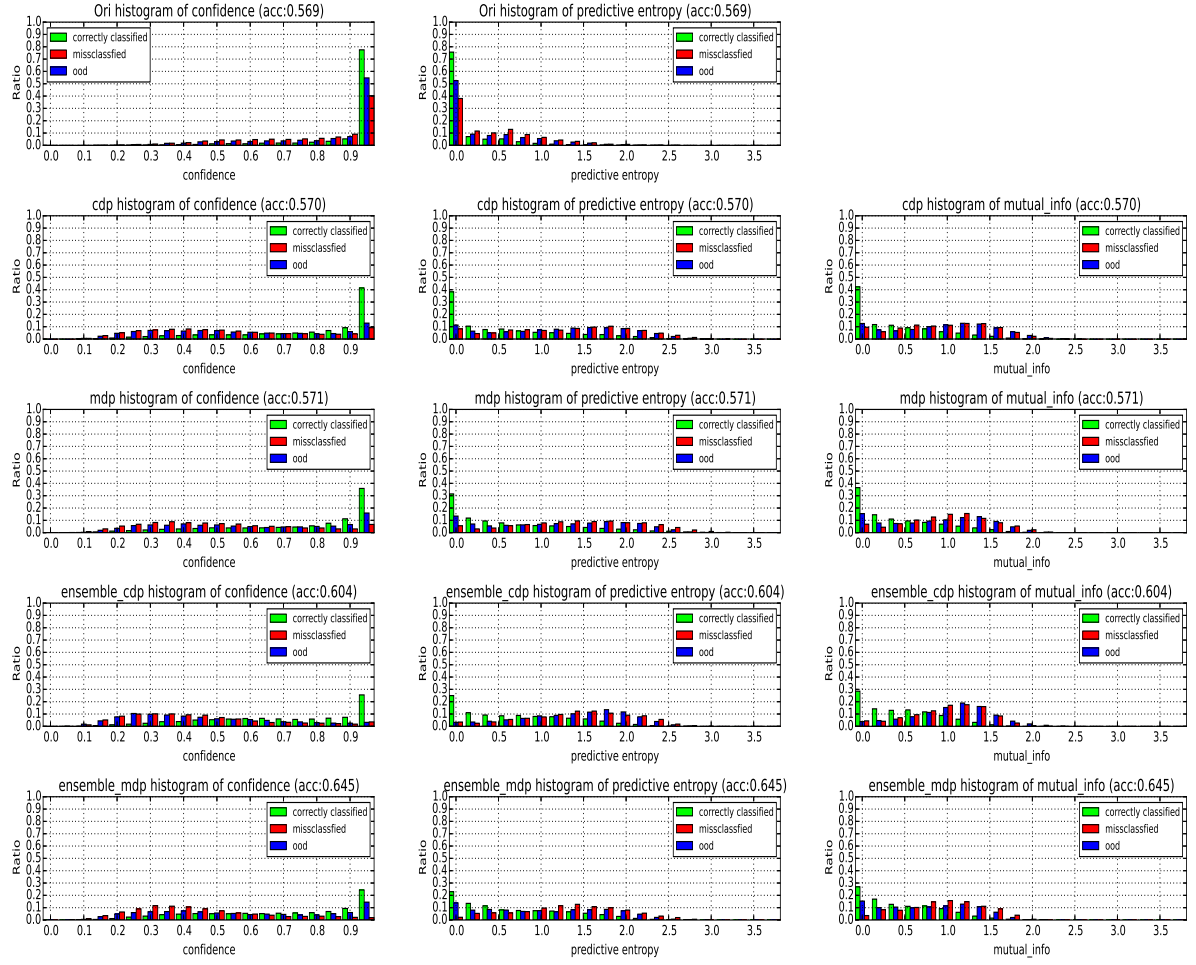


Figure 4.7: Uncertainty histograms of confidence, predictive entropy, mutual information for ori, cdp, mdp, ensemble of cdp, ensemble of mdp from one of three runs.

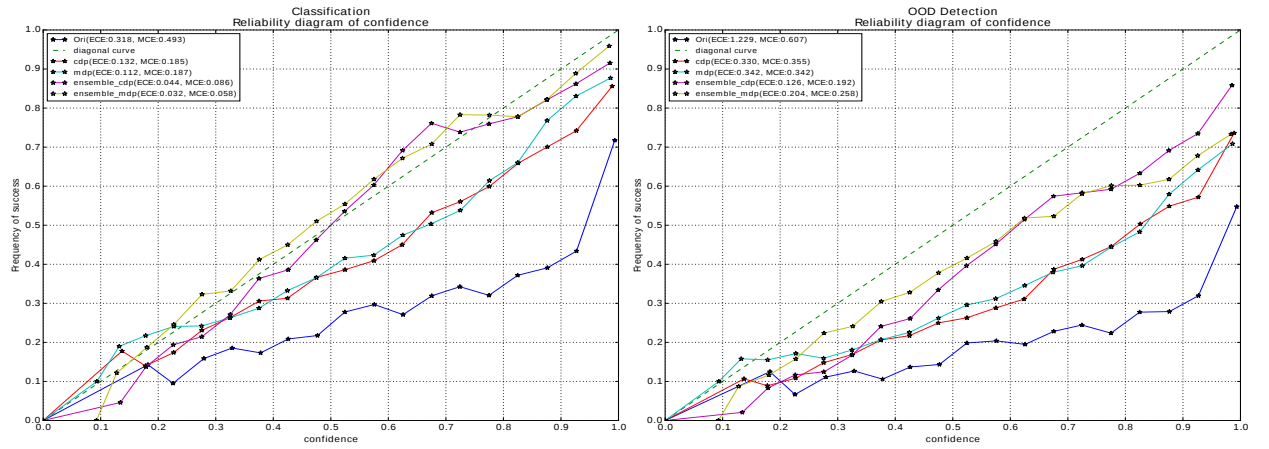


Figure 4.8: Calibration curve for ori, cdp, mdp, ensemble of cdp, ensemble of mdp from one of three runs.

ROC curve and PR curve When it comes to the separability metrics in figure 4.9, as stated before, we measure separability between correct prediction and miss-classification (on the left) and that between correct prediction and OOD prediction (on the right). On separability between correction prediction and miss-classification, cdp, mdp and their ensemble version can improve the results compared with original version in both ROC curve and PR curve on the left. On the other hand, the improvements on separability between correct prediction and OOD data are more obvious on the right. However, as observed before, the ensemble of mdp does not work well with OOD data, which can be observe with this metrics. We can know more from the curves that the worse performance with OOD data in mdp is enlarged by building ensemble of it.

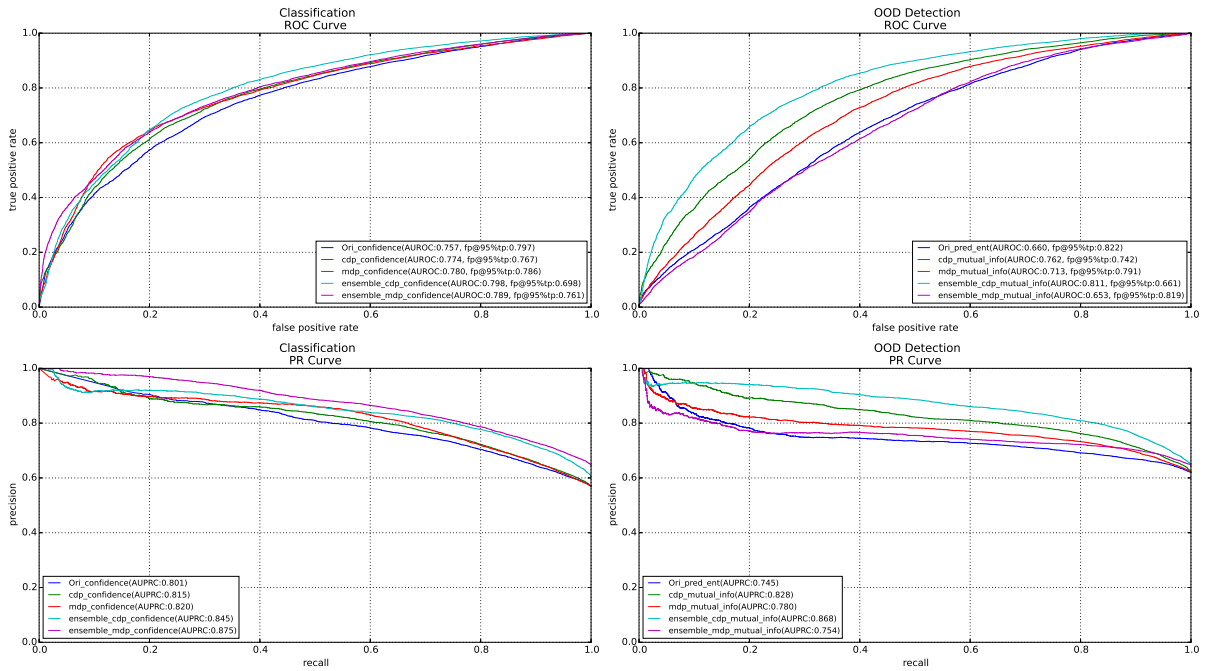


Figure 4.9: ROC and PR curve of ori, cdp, mdp, ensemble of cdp, ensemble of one of three runs.

Before we only show the visual metrics for evaluation, it's better to quantify the performance quantitatively. We adopt different aforementioned metrics to achieve this which are presented in the following table 4.1. In order to exclude the influence of random initialization, we train three models of different random seeds and report the mean and standard deviation for each metric. We can see the improvements by cdp, mdp and their ensemble versions along all the metrics. The phenomenon that mdp does not work well with OOD and this effect enlarged by ensemble is also expressed in the table. Based on these experiment results, we can draw the conclusion that, concrete dropout and multiple dropout can improve both accuracy and calibration performance compared with original ResNet50. While the ensemble of concrete dropout and multiple dropout work even better without OOD data, the ensemble of multiple dropout work worse than that of concrete dropout

when considering OOD data in this work. When comparing concrete dropout and multiple dropout, concrete dropout is more robust than multiple dropout because the latter one does not work well with OOD data.

Table 4.1: Quantitative results of acc, bs, nll, ece, mce, auROC, auPR averaged from 3 different random seeds

	accuracy \uparrow	brier_score \downarrow	negative log likelihood \downarrow
ori	0.568 \pm 0.008	0.722 \pm 0.019	3.242 \pm 0.340
cdp	0.577 \pm 0.008	0.594 \pm 0.013	2.088 \pm 0.181
mdp	0.599 \pm 0.023	0.566 \pm 0.020	1.940 \pm 0.064
emsemble_cdp	0.604	0.534	1.452
emsemble_mdp	0.645	0.496	1.389

	expected calibration error(w/o. OOD/ w. OOD) \downarrow	maximal calibration error(w/o. OOD/ w. OOD) \downarrow	area under ROC (vs. Miss- classified/ vs. OOD) \uparrow	area under PR curve (vs. Miss- classified/ vs. OOD) \uparrow
ori	0.304 \pm 0.016 / 0.633 \pm 0.065	0.461 \pm 0.027 / 0.362 \pm 0.025	0.750 \pm 0.007 / 0.664 \pm 0.011	0.802 \pm 0.008 / 0.751 \pm 0.018
cdp	0.124 \pm 0.023/ 0.288 \pm 0.048	0.206 \pm 0.015/ 0.374 \pm 0.018	0.775 \pm 0.008/ 0.783 \pm 0.022	0.825 \pm 0.007/ 0.850 \pm 0.022
mdp	0.114 \pm 0.012/ 0.383 \pm 0.046	0.199 \pm 0.016/ 0.367 \pm 0.023	0.780 \pm 0.011/ 0.709 \pm 0.004	0.838 \pm 0.013/ 0.788 \pm 0.006
emsemble_cdp	0.044/ 0.042	0.086/ 0.093	0.798/ 0.811	0.845/ 0.868
emsemble_mdp	0.032/ 0.170	0.058/ 0.227	0.789/ 0.653	0.875/ 0.754

Comparison with Laplace approximation

Because Laplace approximation requires only MAP point estimate of model parameter. Therefore we take the already trained model of different approaches as our MAP point estimate and compute the approximation of Kronecker factors with half of training set. We set the scale parameter of Kronecker factors \sqrt{N} as 1 and dump factor $\sqrt{\tau}$ as 15 in equation 2.39 based on grid search on the validation set.

In the following figures, we show the histograms, calibration curves, ROC curve and PR curve of cdp, mdp and their Laplace approximation versions. We can see that the Laplace approximation can achieve similar results as the cdp or mdp do in uncertainty histograms, calibration curves as well as ROC curve and PR curve. The histograms of different uncertainty measure show similar trend as before.

If we compare the Laplace approximation between cdp and mdp, we can see that the result of Laplace approximation is similar to the their dropout versions. This can be observed through that the percentage of OOD prediction with high confidence of Laplace approximation for network trained with mdp is higher than that of Laplace approximation for network trained with cdp in the histogram. Besides, in ROC curve/PR curve on correct prediction vs. OOD data prediction of mdp and cdp, the area under curve of Laplace approximation of mdp is obviously lower than that of cdp. This phenomenon also exists in the histogram and ROC as well as PR curve of mdp and cdp. In the table of quantitative results (cf. 4.2.2), the quantities of different metrics corresponding to the visual metrics. To summarize this, although Laplace can achieve similar results as concrete dropout and multiple dropout do, their performance is highly related to the point estimate of parameter obtained during training.

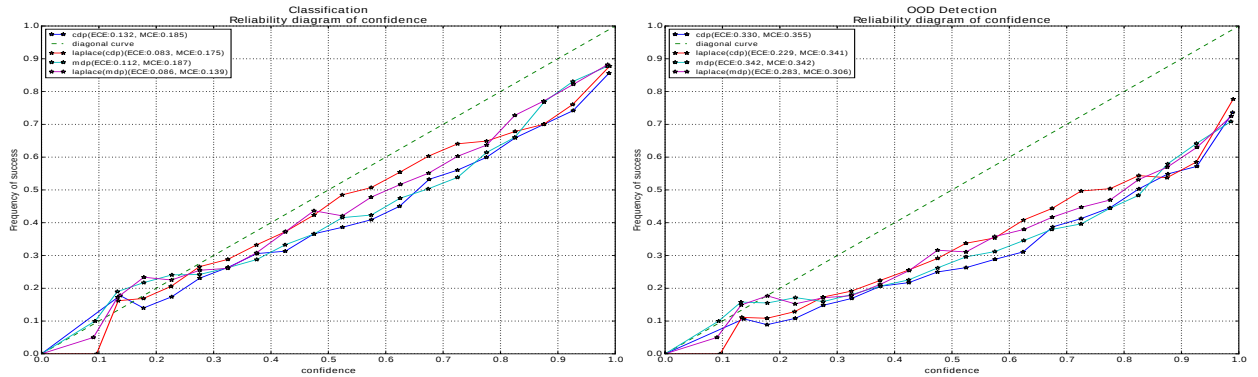


Figure 4.10: Calibration curve of cdp, mdp and their Laplace approximation versions from one of three runs.

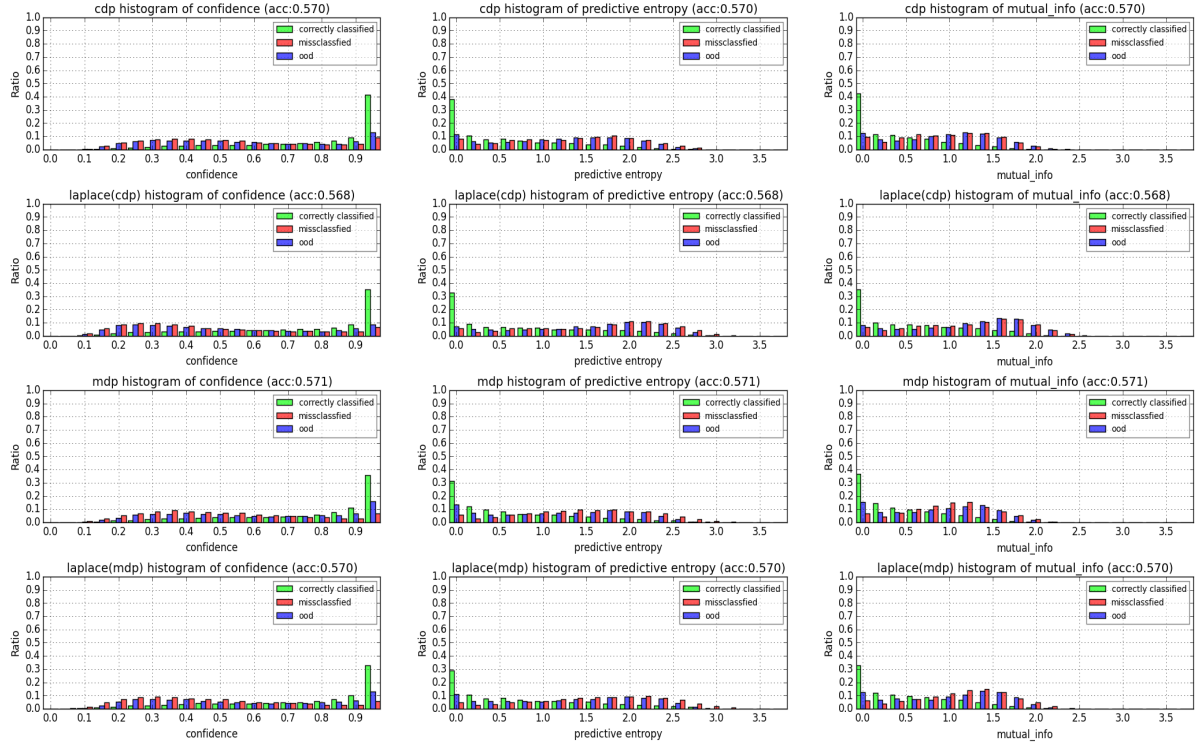


Figure 4.11: Uncertainty histogram of cdp, mdp and their Laplace approximation versions with confidence, predictive entropy and mutual information as uncertainty measure in one of three runs.

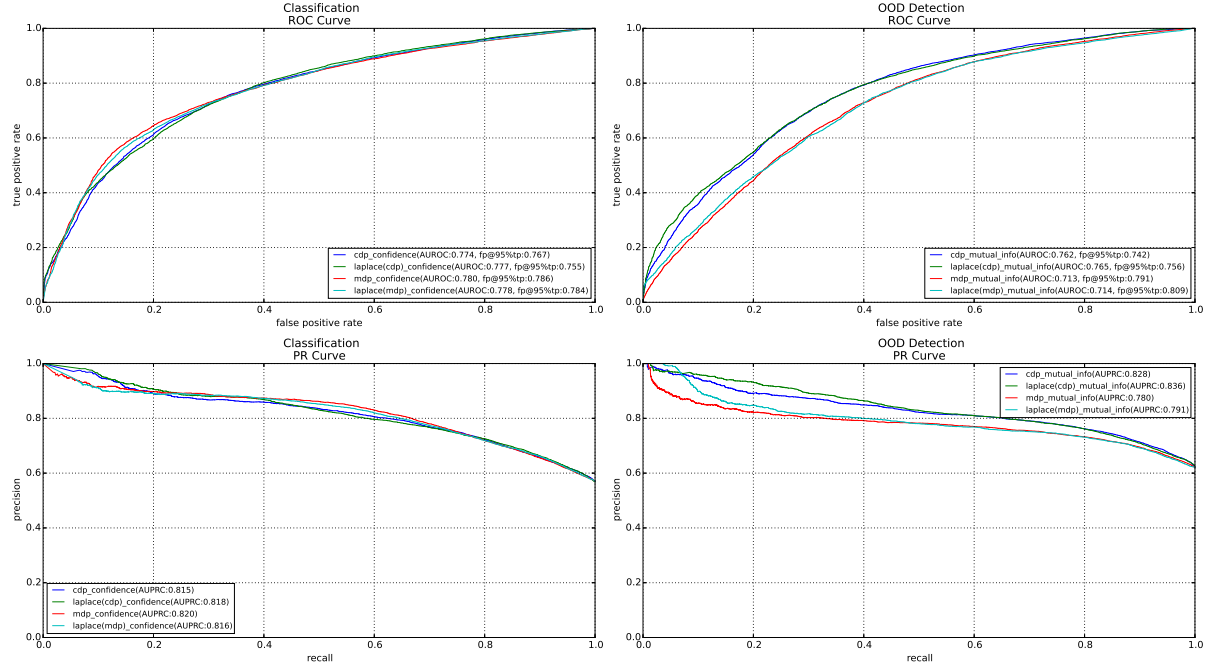


Figure 4.12: ROC and PR curve of cdp, mdp and their Laplace approximation versions from one of three runs.

Table 4.2: Quantitative results of acc, bs, nll, ece, mce, auROC, aupr averaged from 3 different random seeds

	accuracy \uparrow	brier score \downarrow	negative log likelihood \downarrow
cdp	0.577 ± 0.008	0.594 ± 0.013	2.088 ± 0.181
laplace_cdp	0.576 ± 0.009	0.602 ± 0.011	2.322 ± 0.350
mdp	0.599 ± 0.023	0.566 ± 0.020	1.940 ± 0.064
laplace_mdp	0.598 ± 0.024	0.567 ± 0.018	1.970 ± 0.117

	expected calibration error(w/o. OOD/ w. OOD) \downarrow	maximal calibration error(w/o. OOD/ w. OOD) \downarrow	area under ROC (vs. Miss- classified/ vs. OOD) \uparrow	area under PR curve (vs. Miss- classified/ vs. OOD) \uparrow
cdp	$0.124 \pm 0.023 /$ 0.288 ± 0.048	$0.206 \pm 0.015 /$ 0.374 ± 0.018	$0.775 \pm 0.008 /$ 0.783 ± 0.022	$0.825 \pm 0.007 /$ 0.850 ± 0.022
laplace_cdp	$0.129 \pm 0.058 /$ 0.341 ± 0.157	$0.235 \pm 0.073 /$ 0.406 ± 0.070	$0.779 \pm 0.004 /$ 0.782 ± 0.017	$0.826 \pm 0.007 /$ 0.849 ± 0.016
mdp	$0.114 \pm 0.012 /$ 0.383 ± 0.046	$0.199 \pm 0.016 /$ 0.367 ± 0.023	$0.780 \pm 0.011 /$ 0.709 ± 0.004	$0.838 \pm 0.013 /$ 0.788 ± 0.006
laplace_mdp	$0.104 \pm 0.018 /$ 0.352 ± 0.061	$0.179 \pm 0.029 /$ 0.352 ± 0.038	$0.776 \pm 0.012 /$ 0.711 ± 0.005	$0.837 \pm 0.015 /$ 0.798 ± 0.005

Ablation study

As we can see from the previous results, mdp can work better than cdp without considering OOD data. However, it underperforms when OOD data is included. In order to investigate the reason behind this, we have done a ablation study to check the influence of aleatoric uncertainty which comes from mainly the feature extractor part. Considering this, we freeze the parameters of feature extractor and just train our probabilistic classifier which is a three-layer MLP. In the following we show the results of this ablation study with the same metrics used before. It's clear that the accuracy would drop when the feature extractor is frozen. However, we can see that mdp can have better calibration performance than cdp and much better than original ResNet50. Regarding the separability metrics, although mdp still doesn't work better than cdp, the gap between is much smaller than the version without freezing feature extractor. These observations are expressed quantitatively in the table 4.2.2.

Based on these observations, we know that the underperformance of mdp on OOD dataset can attribute to the influence of aleatoric uncertainty. If the feature extractor is not

frozen, then the aleatoric uncertainty is trained to adapt to the dataset. The problem why optimizing dropout rates for each hidden units can change the point estimate of the model parameter to have low aleatoric uncertainty for OOD data arises now. One possible reason for this could be the variance of estimation of derivatives of variational parameters is too large because of increased number of variational parameters[KSW15]. And the larger variances are back propagated to the feature extractor which may make the point estimate deviate from the better local minimum. One possible solution to resolve this is to sample the realization of activations of hidden units independent to each data instance, which may decrease the variance and help to improve the results.

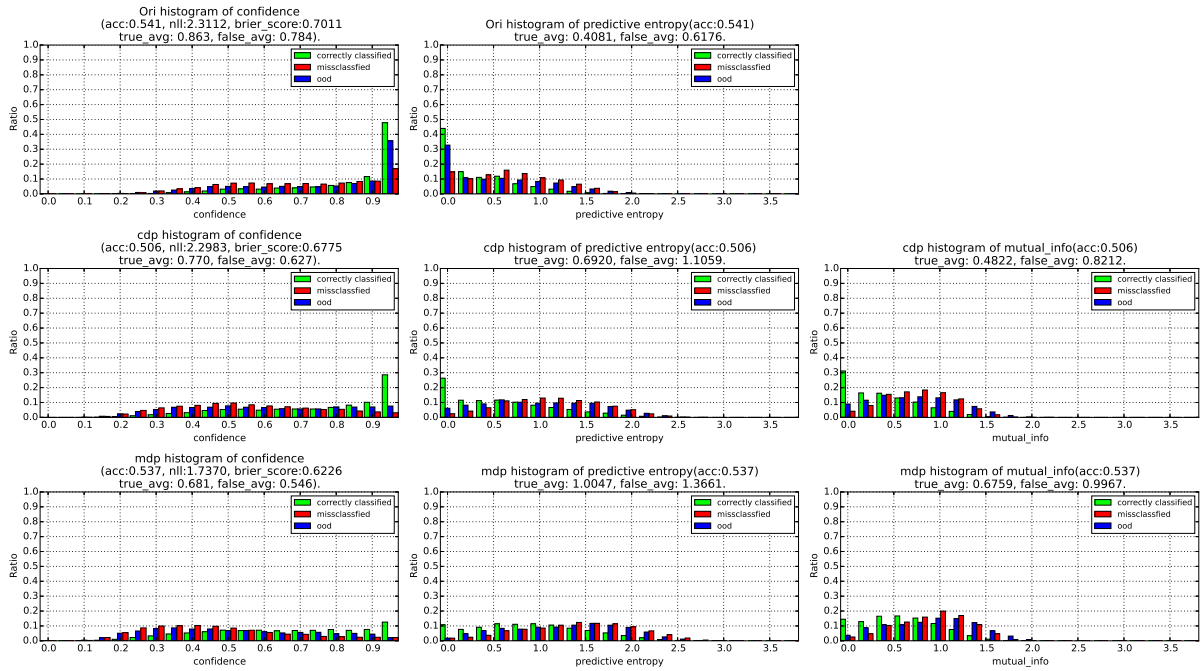


Figure 4.13: Uncertainty histogram of one of three runs.

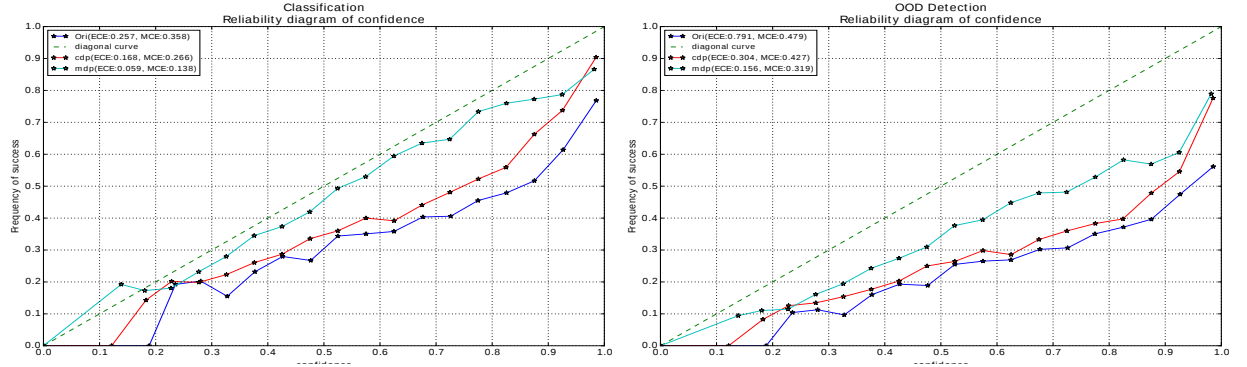


Figure 4.14: Calibration curve of one of three runs.

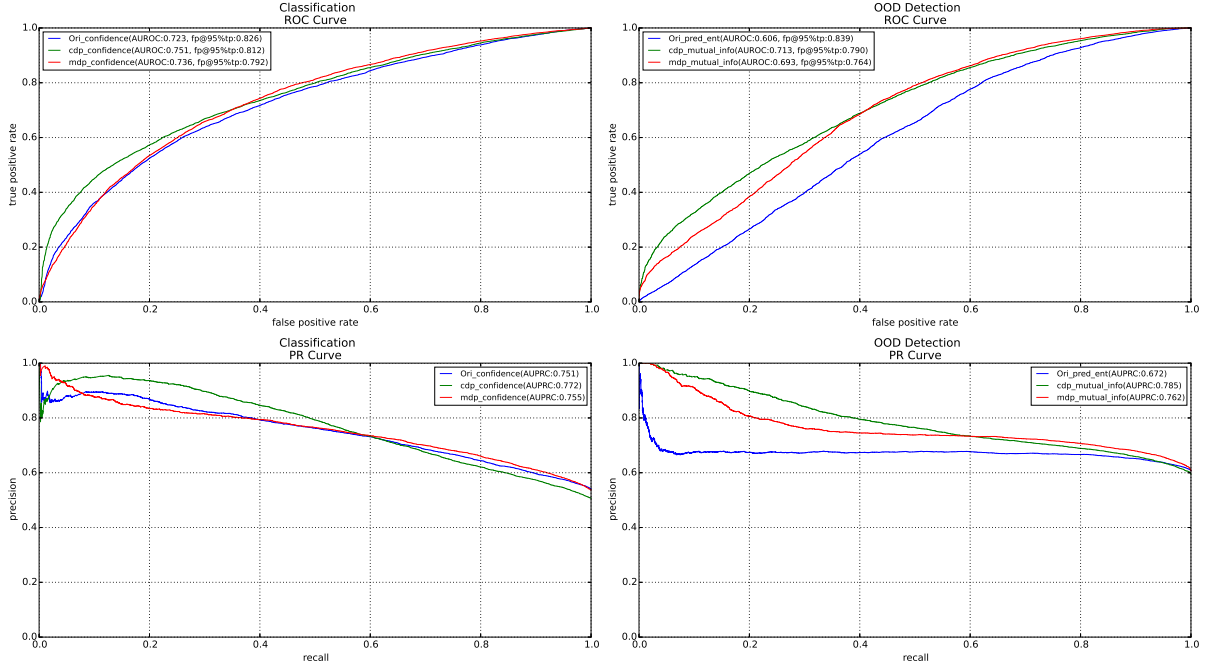


Figure 4.15: ROC and PR curve of one of three runs.

Table 4.3: Quantitative results of acc, bs, nll, ece, mce, auroc, auapr.

	accuracy \uparrow	brier_score \downarrow	negative log likelihood \downarrow
ori	0.532 ± 0.015	0.717 ± 0.023	2.383 ± 0.053
cdp	0.521 ± 0.011	0.655 ± 0.017	2.190 ± 0.149
mdp	0.520 ± 0.012	0.641 ± 0.014	1.804 ± 0.073

	expected calibration error(w/o. OOD/ w. OOD) \downarrow	maximal calibration error(w/o. OOD/ w. OOD) \downarrow	area under ROC (vs. Miss- classified/ vs. OOD) \uparrow	area under PR curve (vs. Miss- classified/ vs. OOD) \uparrow
ori	$0.262 \pm 0.011 /$ 0.793 ± 0.012	$0.365 \pm 0.010 /$ 0.491 ± 0.013	$0.716 \pm 0.005 /$ 0.602 ± 0.014	$0.724 \pm 0.022 /$ 0.681 ± 0.007
cdp	$0.141 \pm 0.022 /$ 0.270 ± 0.037	$0.203 \pm 0.053 /$ 0.359 ± 0.056	$0.748 \pm 0.002 /$ 0.712 ± 0.006	$0.770 \pm 0.002 /$ 0.782 ± 0.004
mdp	$0.068 \pm 0.009 /$ 0.158 ± 0.006	$0.134 \pm 0.018 /$ 0.311 ± 0.008	$0.730 \pm 0.009 /$ 0.674 ± 0.014	$0.749 \pm 0.009 /$ 0.748 ± 0.010

4.3 Automatic labeling experiments

4.3.1 Introduction

In this experiment, our goal is to perform a proof-of-concept experiment, which simulates the situation that a classifier trained on a public large scale dataset or synthetic dataset is deployed in a real world environment and tries to fine-tune itself to improve performance with automatically labeled dataset and as little manually labeled data as possible. The automatically labeled data is collected based on improved uncertainty estimation.

This way to enable the classifier to learn continuously is seldom seen in the literatures. Therefore, in order to investigate the practicability of the way, we firstly perform experiments with restriction on the size of collected dataset.

Then we try to extend this way to another dataset to test its generalization ability in second experiment. We employ data augmentation to address the problem of imbalance in number of data sample of each category found in the first experiment and show the results.

4.3.2 Experiments I

The first experiment is performed on the WRGBD and UniHB dataset. As is mentioned in experiment II in uncertainty estimation, we train classifier with entire WRGBD dataset and evaluate performance of uncertainty estimation on adaptation set from UniHB dataset, which contain the images of object of elevation 30° and 60° , whose size is around 17.1×10^3 . The images of 45° with size 8.6×10^3 are used for final evaluation after fine-tuning the classifier. The dataset used for fine-tuning are collected from the adaptation set with different settings in order to investigate the issues of this way for continuous learning.

In this experiment, the model we use for fine-tuning is classifier trained with concrete dropout. The uncertainty measure we used is confidence because in this case we do not consider OOD data and confidence performs better in separating correct predictions and miss-classifications based on the results of previous experiments.

One major restriction we impose in this experiment is the size of dataset for fine-tuning. Because if we do not restrict the size, the performance is influenced by both the size of dataset and other factors such as imbalance or quality of dataset. In order to exclude the influence of dataset size, we fix it as 3% of the adaptation set which is around 510. If the dataset is balanced in number of data sample of category, then there are around ten data samples in each category. The influence of other factors we want investigate are the imbalance of dataset and quality of data sample in dataset. In the following table, we use capital letter to denote different settings of dataset for fine-tuning.

Here we firstly illustrate the steps of labeling in this experiment which is different to the strategy in experiment II.

For C, D, E, F, G, we conduct the following steps:

1. Automatic labeling: select the most confident predictions and label them with their predictions.
2. Manual labeling: select the data randomly or with least confidence required for manual labeling.
3. Combining: combine automatically labeled and manually labeled data, then check if any category does not have data sample, if yes add one random data sample of this category.
4. Balancing dataset: firstly check if number of data sample of any categories exceed the average number of data sample (which can be calculated beforehand), if yes, drop the data sample of this category randomly until its number reach the average number of category. After that, if augmentation is chosen, then use augmentation to increase of the number of data sample of category in which the number of data sample is less than average number of data sample.

To note that, the accuracy of automatically labeled data in C, E, F, G is around 96%, which means not all labels of them are correct. And "**randomly**" in the settings means that we randomly choose the data according to the number of data sample in each category in order to make the dataset **as balanced as possible**.

Table 4.4: Results of fine-tuning of different settings

Settings of dataset for fine-tuning	Accuracy (average over 3 random seeds)
A: 0% (no fine-tuning)	66.9%
B: 3% manually labeled data randomly (balanced)	91.7%
C: 3% automatically labeled data (imbalanced)	79.0%
D: 3% manually labeled data with least confidence (imbalanced)	83.3%
E: 2% automatically labeled data and 1% manually labeled data with least confidence, augmentation for balance (balanced)	83.8%
F: 2% automatically labeled data and 0.5% manually labeled data with least confidence and 0.5% manually labeled data randomly, augmentation for balance (balanced)	88.3%
G: 2% automatically labeled data and 1% manually labeled data randomly, augmentation for balance (balanced)	89.6%

From the description of settings we know that predictions with most confidence and least confidence are imbalanced in number of data sample of each category. Based on the results in the table above, we can have the following observations:

- diverse and balanced dataset can achieved best performance.
- the difference of performance between C and D may attribute to quality of dataset which include correct labeling and information of data sample.
- by comparing performance of E, F and G, we know that diverse data sample before augmentation is important to achieve better performance. When considering information of data sample, manually labeling of lest confident data sample sounds better. However, we also know that before augmentation dataset in E is more imbalanced than that in F, and F is more imbalanced than G because of the higher proportion of predictions with least confidence. Therefore it shows that data balance is more important than information of data sample in this case.

By comparing the results of other settings with that of A and B, we can see that the manual labeling effort can be reduced based on automatic labeling. The performance of fine-tuned domain specific classifier can nearly reach the the performance of classifier fine-tuned with all manually labeled data.

In the end, we can also draw a conclusion that diversity of data and balance of number of data sample of each category play an significant role. Based on the last aforementioned observation, balance of number of data sample of each category is more important in than information of data sample.

4.3.3 Experiments II

The second experiment is performed on synthetic Tless dataset and original Tless dataset.

4.4 Context-based improvement experiments

4.4.1 Experiments settings

4.4.2 Results

4.4.3 Analysis

Chapter 5

Conclusion

5.0.1 Discussions and conclusions

-

5.0.2 Future work

- check uncertainty distribution of different view points of the object.
- try distillation method to decrease runtime of testing.
- try more complex and expressive approximate posterior distribution such as normalizing flow.
- improve performance of multiple dropout.

Appendix A

Appendix

A.1 KL condition

Proposition 1. Fix $K, L \in \mathbb{N}$, a probability vector $\mathbf{p} = (p_1, \dots, p_L)$, and $\Sigma_i \in \mathbb{R}^{K \times K}$ diagonal positive definite for $i = 1, \dots, L$, with the elements of each Σ_i not dependent on K . Let

$$q(\mathbf{x}) = \sum_{i=1}^L p_i \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \Sigma_i)$$

be a mixture of Gaussians with L components and $\boldsymbol{\mu}_i \in \mathbb{R}^K$, let $p(\mathbf{x}) = \mathcal{N}(0, \mathbf{I}_K)$, and further assume that $\boldsymbol{\mu}_i - \boldsymbol{\mu}_j \sim \mathcal{N}(0, \mathbf{I})$ for all i and j .

The KL divergence between $q(\mathbf{x})$ and $p(\mathbf{x})$ can be approximated as:

$$KL(q(\mathbf{x})||p(\mathbf{x})) \approx \sum_{i=1}^L \frac{p_i}{2} (\boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \text{tr}(\Sigma_i) - K(1 + \log 2\pi) - \log(\det(\Sigma_i))) - \mathcal{H}(\mathbf{p})$$

with

$$\mathcal{H}(\mathbf{p}) = - \sum_{i=1}^L p_i \log p_i$$

for large enough K .

Proof. We have

$$\begin{aligned} KL(q(\mathbf{x})||p(\mathbf{x})) &= \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} \\ &= \int q(\mathbf{x}) \log q(\mathbf{x}) d\mathbf{x} - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} \\ &= -\mathcal{H}(q(\mathbf{x})) - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

which is sum of entropy of $q(\mathbf{x})$ and the expected log probability of \mathbf{x} . The expected log probability can be evaluated analytically, but the entropy term has to be approximated.

We begin by approximating the entropy term. We write

$$\begin{aligned}\mathcal{H}(q(\mathbf{x})) &= - \sum_{i=1}^L p_i \int \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \log q(\mathbf{x}) d\mathbf{x} \\ &= - \sum_{i=1}^L p_i \int \mathcal{N}(\boldsymbol{\epsilon}_i; 0, \mathbf{I}) \log q(\boldsymbol{\mu}_i + \mathbf{L}_i \boldsymbol{\epsilon}_i) d\boldsymbol{\epsilon}_i\end{aligned}$$

with a re-parameterization of \mathbf{x} , that is $\mathbf{x} = \boldsymbol{\mu}_i + \mathbf{L}_i \boldsymbol{\epsilon}_i$ with $\mathbf{L}_i \mathbf{L}_i^T = \boldsymbol{\Sigma}_i$ and $\boldsymbol{\epsilon}_i \sim \mathcal{N}(0, \mathbf{I})$.

Now the term inside logarithm can be written as

$$\begin{aligned}q(\boldsymbol{\mu}_i + \mathbf{L}_i \boldsymbol{\epsilon}_i) &= \sum_{j=1}^L p_j \mathcal{N}(\boldsymbol{\mu}_i + \mathbf{L}_i \boldsymbol{\epsilon}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \\ &= \sum_{j=1}^L p_j (2\pi)^{-\frac{K}{2}} \det(\boldsymbol{\Sigma}_j)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \|\boldsymbol{\mu}_j - \boldsymbol{\mu}_i - \mathbf{L}_i \boldsymbol{\epsilon}_i\|_{\boldsymbol{\Sigma}_j}^2 \right\}\end{aligned}$$

where $\|\cdot\|_{\boldsymbol{\Sigma}}$ is Mahalanobis distance with $\boldsymbol{\Sigma}$ as covariance matrix. Since $\boldsymbol{\mu}_i - \boldsymbol{\mu}_j$ is assumed to be normally distributed, the quantity $\boldsymbol{\mu}_i - \boldsymbol{\mu}_j - \mathbf{L}_i \boldsymbol{\epsilon}_i$ is also normally distributed. Since the expectation of a generalized χ^2 distribution with K degrees of freedom increases with K , we have that $K \gg 0$ implies that $\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j - \mathbf{L}_i \boldsymbol{\epsilon}_i\|_{\boldsymbol{\Sigma}_j}^2 \gg 0$ for $i \neq j$ (since the elements of $\boldsymbol{\Sigma}_j$ does not depend on K). Finally, we have for $i = j$ that $\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j - \mathbf{L}_i \boldsymbol{\epsilon}_i\|_{\boldsymbol{\Sigma}_j}^2 = \boldsymbol{\epsilon}_i^T \mathbf{L}_i^T \mathbf{L}_i^{-T} \mathbf{L}_i^{-1} \mathbf{L}_i \boldsymbol{\epsilon}_i = \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i$. Therefore the last equation can be approximated as

$$q(\boldsymbol{\mu}_i + \mathbf{L}_i \boldsymbol{\epsilon}_i) \approx p_i (2\pi)^{-\frac{K}{2}} \det(\boldsymbol{\Sigma}_i)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i \right\}$$

I.e. in high dimensions the mixture components will not overlap. This gives us

$$\begin{aligned}\mathcal{H}(q(\mathbf{x})) &\approx - \sum_{i=1}^L p_i \int \mathcal{N}(\boldsymbol{\epsilon}_i; 0, \mathbf{I}) \log \left(p_i (2\pi)^{-\frac{K}{2}} \det(\boldsymbol{\Sigma}_i)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i \right\} \right) d\boldsymbol{\epsilon}_i \\ &= - \sum_{i=1}^L p_i \left[\log p_i - \frac{K}{2} \log 2\pi - \frac{1}{2} \log \det(\boldsymbol{\Sigma}_i) - \frac{1}{2} \int \mathcal{N}(\boldsymbol{\epsilon}_i; 0, \mathbf{I}) \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i d\boldsymbol{\epsilon}_i \right] \\ &= \sum_{i=1}^L \frac{p_i}{2} (K \log 2\pi + \log \det(\boldsymbol{\Sigma}_i) + \int \mathcal{N}(\boldsymbol{\epsilon}_i; 0, \mathbf{I}) \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i d\boldsymbol{\epsilon}_i) + \mathcal{H}(\mathbf{p})\end{aligned}$$

Since $\boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i$ distributes according to a χ^2 distribution, its expectation is K , and in the end the entropy term can be approximated as

$$\mathcal{H}(q(\mathbf{x})) \approx \sum_{i=1}^L \frac{p_i}{2} (K(\log 2\pi + 1) + \log \det(\boldsymbol{\Sigma}_i)) + \mathcal{H}(\mathbf{p})$$

Next, we can evaluate the expected log probability term, we get

$$\int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} = \sum_{i=1}^L p_i \int \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \log p(\mathbf{x}) d\mathbf{x}$$

for $p(\mathbf{x}) = \mathcal{N}(0, \mathbf{I}_K)$, it is easy to show that

$$\begin{aligned} \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} &= \sum_{i=1}^L p_i \int \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \log \left[(2\pi)^{-\frac{K}{2}} \det(\mathbf{I}_K)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} \mathbf{x}^T \mathbf{x}\right\} \right] d\mathbf{x} \\ &= \sum_{i=1}^L p_i \left[-\frac{K}{2} \log 2\pi - \frac{1}{2} \int \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \mathbf{x}^T \mathbf{x} d\mathbf{x} \right] \\ &= -\frac{1}{2} \sum_{i=1}^L p_i (\boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \text{tr}(\boldsymbol{\Sigma}_i) + K \log 2\pi) \end{aligned}$$

Finally, combining the equations above, we have

$$KL(q(\mathbf{x}) || p(\mathbf{x})) \approx \sum_{i=1}^L \frac{p_i}{2} (\boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \text{tr}(\boldsymbol{\Sigma}_i) - K - \log(\det(\boldsymbol{\Sigma}_i))) - \mathcal{H}(\mathbf{p})$$

as required to show. □

A.2 Figures

Beispiel für eine Tabelle:

Table A.1: Beispiel für eine Beschriftung. Tabellenbeschriftungen sind üblicherweise über der Tabelle platziert.

left	center	right
entry	entry	entry
entry	entry	entry
entry	entry	entry

A.3 Implementation Details

List of Figures

2.1	Difference between parameter estimation of deterministic neural network and Bayesian neural network.	10
2.2	How dropout works[SHK ⁺ 14].	13
2.3	A two layer neural network example of dropout inference, a Bernoulli random variable is imposed on each unit of input layer and hidden layer. . . .	14
2.4	Example of Gumbel-max trick for drawing samples from a discrete distribution whose random variable has 3 states and $\{\alpha_i\}_{i=1,2,3}$ as class parameters representing the possibility of occurrence of that class. $\{G_i\}_{i=1,2,3}$ are i.i.d Gumbel(0, 1) [MMT16].	22
2.5	Example of continuous approximation to Gumbel-max trick for drawing samples from a discrete distribution whose random variable has 3 states and $\{\alpha_i\}_{i=1,2,3}$ as class parameters representing the possibility of occurrence of that class. $\{G_i\}_{i=1,2,3}$ are i.i.d Gumbel(0, 1)[MMT16].	23
2.6	One sample and average value of 100 samples drawn from continuous approximation of Bernoulli distribution with parameter $p = [0.1, 0.2, \dots, 1.0]$ and temperature $\lambda = 0.1$	24
2.7	Changes along epochs of keep probability in training network with concrete dropout layer.	26
2.8	Different dropout rates for different hidden units in multi-drop.	26
2.9	Modified network architecture.	28
3.1	Simple example of conditional random field.	36
4.1	Example of masked images of objects from 51 categories in WRGBD and UniHB dataset. In each category, the left is from WRGBD and the right is from UniHB. We randomly pick one instance for the objects of WRGBD. We can see some light and appearance difference between objects in these two datasets.	42
4.2	Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples. . .	42
4.3	Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples. . .	43

4.4	Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples. .	43
4.5	Uncertainty(confidence, predictive entropy, mutual information) histograms of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout.	49
4.6	Calibration curves of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout.	49
4.7	Uncertainty histograms of confidence, predictive entropy, mutual information for ori, cdp, mdp, ensemble of cdp, ensemble of mdp from one of three runs.	52
4.8	Calibration curve for ori, cdp, mdp, ensemble of cdp, ensemble of mdp from one of three runs.	52
4.9	ROC and PR curve of ori, cdp, mdp, ensemble of cdp, ensemble of one of three runs.	53
4.10	Calibration curve of cdp, mdp and their Laplace approximation versions from one of three runs.	55
4.11	Uncertainty histogram of cdp, mdp and their Laplace approximation versions with confidence, predictive entropy and mutual information as uncertainty measure in one of three runs.	56
4.12	ROC and PR curve of cdp, mdp and their Laplace approximation versions from one of three runs.	56
4.13	Uncertainty histogram of one of three runs.	58
4.14	Calibration curve of one of three runs.	59
4.15	ROC and PR curve of one of three runs.	59

List of Tables

4.1	Quantitative results of acc, bs, nll, ece, mce, auROC, aupr averaged from 3 different random seeds	54
4.2	Quantitative results of acc, bs, nll, ece, mce, auROC, aupr averaged from 3 different random seeds	57
4.3	Quantitative results of acc, bs, nll, ece, mce, auROC, aupr.	59
4.4	Results of fine-tuning of different settings	62
A.1	Beispiel für eine Beschriftung. Tabellenbeschriftungen sind üblicherweise über der Tabelle platziert.	67

Bibliography

- [ABC⁺16] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.
- [AOS⁺16] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [BCKW15] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- [Bis06] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [BRMW15] Anoop Korattikara Balan, Vivek Rathod, Kevin P Murphy, and Max Welling. Bayesian dark knowledge. In *Advances in Neural Information Processing Systems*, pages 3438–3446, 2015.
- [DKD09] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural Safety*, 31(2):105–112, 2009.
- [DL91] John S Denker and Yann Lecun. Transforming neural-net output levels to probability distributions. In *Advances in neural information processing systems*, pages 853–859, 1991.
- [DT18] Terrance DeVries and Graham W Taylor. Learning confidence for out-of-distribution detection in neural networks. *arXiv preprint arXiv:1802.04865*, 2018.
- [FCSG17] Reuben Feinman, Ryan R Curtin, Saurabh Shintre, and Andrew B Gardner. Detecting adversarial samples from artifacts. *arXiv preprint arXiv:1703.00410*, 2017.

- [FRD18] Di Feng, Lars Rosenbaum, and Klaus Dietmayer. Towards safe autonomous driving: Capture uncertainty in the deep neural network for lidar 3d vehicle detection. *arXiv preprint arXiv:1804.05132*, 2018.
- [Gal] Yarin Gal. *Uncertainty in deep learning*. PhD thesis.
- [GBR07] Tilman Gneiting, Fadoua Balabdaoui, and Adrian E Raftery. Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2):243–268, 2007.
- [GG16] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.
- [GHK17] Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. In *Advances in Neural Information Processing Systems*, pages 3581–3590, 2017.
- [GIG17] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. *arXiv preprint arXiv:1703.02910*, 2017.
- [GMR] Yarin Gal, Rowan McAllister, and Carl Edward Rasmussen. Improving pilco with bayesian neural network dynamics models.
- [GN99] AK Gupta and DK Nagar. *Matrix Variate Distributions*, volume 104. CRC Press, 1999.
- [GPSW17] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. *arXiv preprint arXiv:1706.04599*, 2017.
- [GR07] Tilman Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [Gra11] Alex Graves. Practical variational inference for neural networks. In *Advances in neural information processing systems*, pages 2348–2356, 2011.
- [GRB08] Carolina Galleguillos, Andrew Rabinovich, and Serge Belongie. Object categorization using co-occurrence, location and appearance. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.
- [HC68] John M Hammersley and Peter Clifford. Markov fields on finite graphs and lattices. 1968.
- [HG16] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- [HHGL11] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel.

- Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- [HHO⁺17] Tomáš Hodaň, Pavel Haluza, Štěpán Obdržálek, Jiří Matas, Manolis Lourakis, and Xenophon Zabulis. T-LESS: An RGB-D dataset for 6D pose estimation of texture-less objects. *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2017.
- [HLA15] José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, pages 1861–1869, 2015.
- [HVC93] Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pages 5–13. ACM, 1993.
- [HZRS16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [JGP16] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [KC16] Alex Kendall and Roberto Cipolla. Modelling uncertainty in deep learning for camera relocalization. In *2016 IEEE international conference on Robotics and Automation (ICRA)*, pages 4762–4769. IEEE, 2016.
- [KFB09] Daphne Koller, Nir Friedman, and Francis Bach. *Probabilistic graphical models: principles and techniques*. 2009.
- [KG17] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017.
- [KGB16] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.
- [KGC] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics.
- [KK11] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, pages 109–117, 2011.
- [KSW15] Durk P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems*, pages 2575–2583, 2015.

- [KW13] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [LAA⁺17] Christian Leibig, Vaneeda Allken, Murat Seçkin Ayhan, Philipp Berens, and Siegfried Wahl. Leveraging uncertainty information from deep neural networks for disease detection. *Scientific reports*, 7(1):17816, 2017.
- [LBRF11] Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *2011 IEEE international conference on robotics and automation*, pages 1817–1824. IEEE, 2011.
- [LLLS17] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *arXiv preprint arXiv:1711.09325*, 2017.
- [LLS17] Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- [LMP01] John Lafferty, Andrew McCallum, and Fernando CN Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. 2001.
- [LPB17] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pages 6402–6413, 2017.
- [LRKT10] Lubor Ladicky, Chris Russell, Pushmeet Kohli, and Philip HS Torr. Graph cut based inference with co-occurrence statistics. In *European Conference on Computer Vision*, pages 239–253. Springer, 2010.
- [LSVDHR16] Guosheng Lin, Chunhua Shen, Anton Van Den Hengel, and Ian Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3194–3203, 2016.
- [LW16] Christos Louizos and Max Welling. Structured and efficient variational deep learning with matrix gaussian posteriors. In *International Conference on Machine Learning*, pages 1708–1716, 2016.
- [LW17] Christos Louizos and Max Welling. Multiplicative normalizing flows for variational bayesian neural networks. *arXiv preprint arXiv:1703.01961*, 2017.
- [Mac92] David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- [MG15] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417, 2015.

- [MGK⁺17] Rowan McAllister, Yarin Gal, Alex Kendall, Mark Van Der Wilk, Amar Shah, Roberto Cipolla, and Adrian Vivian Weller. Concrete problems for autonomous vehicle safety: Advantages of bayesian deep learning. International Joint Conferences on Artificial Intelligence, Inc., 2017.
- [Min01] Thomas P Minka. Expectation propagation for approximate bayesian inference. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 362–369. Morgan Kaufmann Publishers Inc., 2001.
- [MMT16] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- [MTM14] Chris J Maddison, Daniel Tarlow, and Tom Minka. A* sampling. In *Advances in Neural Information Processing Systems*, pages 3086–3094, 2014.
- [Nea12] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [OBPVR16] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped dqn. In *Advances in neural information processing systems*, pages 4026–4034, 2016.
- [Pea14] Judea Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.
- [QCRS⁺05] Joaquin Quinonero-Candela, Carl Edward Rasmussen, Fabian Sinz, Olivier Bousquet, and Bernhard Schölkopf. Evaluating predictive uncertainty challenge. In *Machine Learning Challenges Workshop*, pages 1–27. Springer, 2005.
- [RBB18] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural networks. 2018.
- [RSGGJ15] J. R. Ruiz-Sarmiento, Cipriano Galindo, and Javier González-Jiménez. Up-gmpp: a software library for contextual object recognition. In *3rd. Workshop on Recognition and Action for Scene Understanding (REACTS)*, 2015.
- [RV09] Nikhil Rasiwasia and Nuno Vasconcelos. Holistic context modeling using semantic co-occurrences. 2009.
- [RVG⁺07] Andrew Rabinovich, Andrea Vedaldi, Carolina Galleguillos, Eric Wiewiora, and Serge Belongie. Objects in context. In *Computer vision, 2007. ICCV 2007. IEEE 11th international conference on*, pages 1–8. IEEE, 2007.
- [SBD⁺14] Robin Senge, Stefan Bösner, Krzysztof Dembczyński, Jörg Haasenritter, Oliver Hirsch, Norbert Donner-Banzhoff, and Eyke Hüllermeier. Reliable classification: Learning classifiers that distinguish aleatoric and epistemic uncertainty. *Information Sciences*, 255:16–29, 2014.

- [SCC17] Shengyang Sun, Changyou Chen, and Lawrence Carin. Learning structured weight uncertainty in bayesian neural networks. In *Artificial Intelligence and Statistics*, pages 1283–1292, 2017.
- [SG18] Lewis Smith and Yarin Gal. Understanding measures of uncertainty for adversarial example detection. *arXiv preprint arXiv:1803.08533*, 2018.
- [Sha48] Claude Elwood Shannon. A mathematical theory of communication. *Bell system technical journal*, 27(3):379–423, 1948.
- [SHK⁺14] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [SM⁺12] Charles Sutton, Andrew McCallum, et al. An introduction to conditional random fields. *Foundations and Trends® in Machine Learning*, 4(4):267–373, 2012.
- [Wil92] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [WT11] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.
- [WVL⁺18] Kuan-Chieh Wang, Paul Vicol, James Lucas, Li Gu, Roger Grosse, and Richard Zemel. Adversarial distillation of bayesian neural network posteriors. *arXiv preprint arXiv:1806.10317*, 2018.
- [ZSDG17] Guodong Zhang, Shengyang Sun, David Duvenaud, and Roger Grosse. Noisy natural gradient as variational inference. *arXiv preprint arXiv:1712.02390*, 2017.