



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(DLR), Maximilian Durner(DLR)
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 23, 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 23, 2019

Place, Date

Signature

Abstract

Deep learning based classifiers have achieved tremendous success on different tasks. However, this kind of classifier cannot provide reliable uncertainty estimation about their predictions and express excessive overconfidence, which means that the models are highly certain although they are not familiar with the related input. This behavior can easily lead to severe consequences in safety-critical applications such as diagnosis of diseases, perception of self-driving car and so on. Especially in the case of the deployment of robots, the robots should be aware of the correctness of the predictions they make, so that they are able to avoid unnecessary accidents and adapt itself in an unfamiliar environment. Besides uncertainty for independent data examples, uncertainty related to dependencies between data examples should be taken into account. When recognizing objects with similar appearances but different contexts in a specific scene, the classifier should express reliably high uncertainty about its predictions because of appearance ambiguities. At the same time, contextual information should be handled properly and utilized to resolve this kind of ambiguities, aiming to further improve the performance.

In this work, Bayesian neural networks with different inference techniques such as dropout variational inference and scalable Laplace approximation are employed to improve the uncertainty estimation. With a reliable uncertainty estimation, a classifier trained on an easily obtainable dataset should be able to adapt itself in a test environment by collecting dataset for fine-tuning with as little manual efforts as possible. Additionally, a conditional random field is employed to further improve the performance by utilizing more information brought by better uncertainty estimation and handling the contextual information between objects properly. Extensive experiments are performed to show that uncertainty estimation can be improved with Bayesian neural networks in terms of different evaluation metrics. Besides, it is demonstrated that manual efforts in collecting dataset for fine-tuning a classifier can be reduced with the help of improved uncertainty estimation on experiments with two different datasets. Last but not least, experiments also show that the conditional random field is able to further improve the performance by making use of better uncertainty estimation and contextual information.

List of Abbreviations

BNN	Bayesian Neural Network.
CRF	Conditional Random Field.
DNN	Deep Neural Network.
KL-div	Kullbach-Leibler divergence.
MAP	Maximum a posterior.
MCMC	Markov Chain Monte Carlo.
MLE	Maximum Likelihood Estimation.
OOD	Out-of-distribution.
SGD	Stochastic Gradient Descent.
SGLD	Stochastic Gradient Langevin Dynamics.
VI	Variational Inference.

Contents

Contents	iii
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?	3
1.4 Contributions and Structure	5
2 Theoretical Background	7
2.1 Bayesian neural network	7
2.1.1 Introduction	7
2.1.2 Dropout variational inference	9
2.1.3 Laplace approximation	23
2.2 Conditional random field	27
2.2.1 Definition	28
2.2.2 Learning	30
2.2.3 Inference	32
3 Technical Approach	33
3.1 Multiple dropout	33
3.2 Modified network architecture	35
3.3 Combination with CRF	37
3.4 Approach for continuous learning	38
4 Experiment and Discussion	40
4.1 Background	40
4.1.1 Dataset	41
4.1.2 Uncertainty measure	45
4.1.3 Evaluation metric	46
4.2 Uncertainty estimation experiments	49
4.2.1 Experiments I: uncertainty estimation on instance recognition . . .	49
4.2.2 Experiments II: uncertainty estimation on category recognition . . .	51
4.3 Automatic labeling experiments	61

4.3.1	Experiment I: evaluation on WRGBD and UniHB dataset	62
4.3.2	Experiment II: evaluation on T-LESS dataset	64
4.4	Context-based improvement experiments	68
4.4.1	Experiment I: evaluation on subset of T-LESS dataset	69
4.4.2	Experiment II: evaluation on entire T-LESS dataset	70
5	Summary and Conclusion	71
5.1	Summary	71
5.2	Conclusion	73
5.3	Future work	73
A	Appendix	75
A.1	Fisher information	75
A.2	Derivation of Re-parameterization Trick	76
A.3	Derivation of Hessian of Log Posterior	77
	List of Figures	78
	List of Tables	81
	Bibliography	82

Chapter 1

Introduction

1.1 What is uncertainty?

This question seems a little philosophical and difficult to answer. However, in the context of probabilistic modeling in machine learning, one reasonable and illustrative answer can be found in [KFB09]: "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". In spite of the innate nondeterminism, the uncertainty can be further categorized into two main types: the **epistemic uncertainty** and the **aleatoric uncertainty** [DKD09, SBD⁺14, KG17].

The former one, epistemic or model uncertainty illustrates the the uncertainty in modeling, associated with **imperfect models** of the real world due to the insufficient or imperfect knowledge of reality, i.e., it refers to the epistemic state of the decision maker. Therefore it can reduced by collecting more knowledge. One simple example to illustrate this kind of uncertainty is a linear regression with a limited number of observed points, where different curves are obtained during training and then used to predict the unobserved points. After observing enough points in some intervals, the generated models or curves are more restricted, hence the uncertainty decreases. On the other hand, if there are insufficient points in specific intervals, upon which the large number of curves can be obtained to fit them during training, which induces high model uncertainty. In reality, it is common to have inadequate amount of data, and thus hard to model the underlying distribution well and precisely. Therefore the observed data can be explained by a large variety of different models, from where the model uncertainty comes.

The latter one, aleatoric or data uncertainty, accounts for the **inherent randomness** of an underlying phenomenon which is expressed as variability in the observed data and meanwhile associated with the limitation of our observing ability. It is treated as non-reducible within our ability to observe the data. Examples for this kind of uncertainty are the noise induced by limited resolution of sensors or the ambiguity from the content of the

data which are unrelated to the model used to represent the data. One simple example can be the classification of images captured by cameras with different resolutions. Images with lower resolutions have the higher uncertainty because they are more ambiguous and less illustrative compared to those with higher resolutions.

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

1.2 Why do we need uncertainty?

The aforementioned model uncertainty represents the belief of a model about its predictions, which brings additional information about the randomness induced by the model compared with the deterministic counter part, which provides only a hard prediction. This kind of uncertainty plays an important role in the two following scenarios.

The first use-cases are those requiring **high safety guarantee** or **optimal decision-making**, in which false predictions can cause hazardous consequences, which is the so-called AI safety[AOS⁺16]. A reliable model uncertainty estimation can equip the model with the ability to know when it does not know. More than that, people or operator can take corresponding counter measures to avoid unnecessary accidents. With more and more machine learning algorithms applied in real life applications, this kind of ability meanwhile plays an more and more important role. Concrete examples include the steering control in self-driving cars[MGK⁺17], disease diagnosis in the medical domain[LAA⁺17], or even applications in aerospace or other domains requiring higher precision and stronger robustness. One could argue that if the trained model can perform perfectly, thus having achieved 100% accuracy, there is no need for the uncertainty information. [DL91] argue that, it is hard or even impossible to have enough training data to define a precise model, which is in coincidence with reality because it is hard to define a task very precisely and get access to enough samples drawn from its real distribution. The uncertainty information is even more valuable and important when algorithms are confronted with multiple tasks[KGC]. Other relevant tasks such as misclassified and Out-of-distribution (OOD) data detection[HG16] and adversarial attacks [KGB16, FCSG17] have also attracted more and more research interests in uncertainty estimation, aiming to address these challenges.

The second scenario are cases requiring **interaction between algorithms and people or environments**. The model uncertainty can be treated as a bridge between machine learning algorithms and humans. Such interaction or the exploration of the environment can establish a more robust and more data-efficient machine learning system. This scenario involves different kind of fields such as active learning[GIG17], reinforcement learning[BCKW15][OBPVR16][GMR], automatic labeling, industrial components inspection and so on. The idea behind that is, with reliable model uncertainty estimation we

know which data sample the model is unfamiliar with and thus providing more information for the model. In active or reinforcement learning, these unfamiliar data samples can be used to improve the corresponding model more efficiently. In tasks involving human robot interaction, the uncertainty estimation indicates when the model requires help from human experts. For instance, in industrial component inspection task, firstly a model which can perform quite well on some evaluation datasets is trained. With the assumption that all predictions made by the model are correct, the system is deployed to the real application. However, there is no guarantee that this model represents the real distribution of data perfectly and is robust against a slight or even large domain gap. If the distributions of the real world data vary a lot because of unexpected factors e.g. equipment aging or changes of weather condition, the model could 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 assists 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 a 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 a semantic view point, which we can be also incorporated into training and yield more robust models. This kind of uncertainty is not quantified explicitly but can be incorporated conceptually with the help of specifically designed models and high level statistic cues or features to **disambiguate predictions** directly. One example of this is the modeling contextual information among pixels by a so-called Conditional Random Field(CRF) in the semantic segmentation task[KK11][SM⁺12][LSVDHR16]. The idea behind that is simple and CRF provides a principle way to achieve this idea. Since uncertainty caused by the appearance ambiguity represents strong correlation between different categories, CRF can model the correlations or dependencies between random variables by taking into account high order potentials. In the task of semantic segmentation, for 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 already known, deep learning is a powerful black box model while it has achieved tremendous success in different tasks. Therefore to obtain and handle uncertainty can help to get a better understanding about this black box model and make deep learning based applications 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, due to the intractability of the real model uncertainty in Deep Neural Network (DNN), there are some **ad-hoc approaches** that try to obtain such a quality measure on specific tasks. [LLS17] combine temperature scaling and input preprocessing which is called adversarial training together and yields a simple but effective OOD sample detector, but this approach requires an additional OOD dataset for training which is hard to collect in reality. [DT18] modify the loss function to train a classifier against OOD data, however this approach is hard to generalize in case of slight domain transfers. [LLS17] treat a uniform distribution as ground truth distribution of the OOD data and show that OOD data can be generated via a so-called Generative Adversarial Network(GAN). This approach seems practical, but training a GAN effectively does not always work for images because of the high dimensionality of the data distribution.

On the other hand, there are some other more theoretically sound approaches which employ Bayesian Neural Network (BNN)[Mac92][Nea12], bootstrapping[OBPVR16], ensemble methods[LPB17] or even ensemble of [SG18]. Although a BNN provides a principal way to obtain model uncertainty by considering distribution on model parameters, it is difficult to infer the exact posterior distribution over model parameters of DNN because of its complex architecture and large scale training set. Therefore approximate inference plays an important role in addressing this issue. The golden standard for the 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(SGLD)[WT11], the latter one attracts a lot of research interests like [BRMW15][WVL⁺18]. However, since it is 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 (KL-div) between the approximate posterior distribution and real posterior. Therefore approximate distribution family needs to be chosen and thus has a large influence on the performance. [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 "re-parameterization trick"[KW13]. HLA15 also used the same approximate distribution but with expectation propagation[Min01] instead of VI. Since the assumption of fully factorization may impose an strong restriction on the flexibility of the approximate distribution and the correlations between weights can not be captured. 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 with non-diagonal covariance matrix [KSW15] which can capture variances of rows in weight matrices. Another more expressive one is matrix variate Gaussian distribution[LW16][SCC17][ZSDG17], which capture both rows and columns correlations in weight matrices. 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 with help of "local re-parameterization trick" [KSW15].

Another alternative is the Laplace approximation [Mac92], which puts a Gaussian centered at Maximum a posterior (MAP) estimate of the model parameters with a covariance determined by the inverse of the Hessian matrix of the log posterior distribution. The main issue of this approach for large neural networks is infeasible computation and storage caused by large size of the Hessian matrix which is quadratic of the number of parameters and inversion operation on it whose complexity is cubic over the number of elements of Hessian matrix. Recently [RBB18] address this problem by approximating the Hessian or the Fisher matrix of log-likelihood with kronecked-factorization approximation which reduces the storage size and complexity of inversion operation. The resulting posterior can also be treated as a 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 many practical applications. Therefore the so-called dropout inference [GG16] looks promising considering this restriction and thus is adopted to obtain model uncertainty in this work. On the other hand, scalable Laplace approximation [RBB18] has the potential to resolve 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 a deep neural network training. Therefore another model, **CRF** [LMP01], is adopted based on the predictions of a discriminative neural network classifier. As mentioned in the previous section, the choice of features of second order potentials is a key factor in this problem. In this work, addressing the scene object recognition task, the correlation between objects in specific scene are of importance. The correlation can be represented by a co-occurrence matrix [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 BNN to obtain reliable and calibrated uncertainty. In order to keep the ability of powerful feature extraction of DNN, we modify a classifier with ResNet50 [HZRS16] as backbone to incorporate different aforementioned inference techniques including dropout variational inference and scalable Laplace approximation. We evaluate prediction and calibration performance of BNN and then select predictions with low uncertainty/high confidence as automatically labeled data based on improved uncer-

tainty estimation. This labeled dataset collected with little manual efforts is used for training a more accurate classifier while the performance of original one drops a lot in case of slight domain transfer caused by possibly by changes of light condition, changes of appearance and even changes of recording equipment. This is proof-of-concept idea which tries to prove that improved uncertainty estimation can assist in automatic labeling and reducing manual effort in fine-tuning a more accurate classifier in case of slight domain transfer during deployment of robots.

The second part is to handle the uncertainty caused by appearance ambiguity properly. 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 this end, CRF is selected to capture the contextual relationship between objects with similar appearance but different contexts and at the same time utilize the information brought by better uncertainty estimation obtained in the first part. Key factor to make CRF to work is to design informative pairwise features which should not contradictory to unary feature. In this work, binary co-occurrence matrix of categories as pairwise features are employed.

The main contributions of this work are:

- to show that dropout variational inference and scalable Laplace approximation can obtain better uncertainty information from a large and complex network on a multi-class classification problem with 51 classes, which can be further improved by learning dropout rates during training and their ensemble.
- to show that manual effort in data collection in a slight domain transfer learning task can be reduced by making use of reliable uncertainty estimation. With help of this, classifier can learn continuously in test environment which has domain gap to the training set.
- to show that CRF can improve the results from discriminative classifier by handling uncertainty caused by appearance ambiguity and utilizing information brought by better uncertainty estimation.

The structure of this work is as follows. In order to give the reader a clear theoretical background of techniques applied in this work, the basic theory and ideas of Bayesian neural network and inference techniques such as dropout variational inference and scalable Laplace approximation, and the theory of CRF are briefly reviewed in chapter 2. In chapter 3, technical approaches proposed and employed in this work, including proposed variants of concrete dropout, modified network architecture, combination of BNN and CRF as well as pipeline of continuous learning are introduced. In chapter 4, extensive experiments are performed to evaluate performance of uncertainty estimation of BNN with different inference techniques. In the following, experiments on performance of continuous learning and improvements brought by CRF are performed. Then summary of this work, conclusion and possible future work are presented in chapter 5.

Chapter 2

Theoretical Background

2.1 Bayesian neural network

In this chapter, a brief review on theory of BNN is given. Since calculating the exact posterior distribution over the parameters of large and complex networks is intractable, proper techniques to perform approximate inference are required. Those techniques can be categorized into two main mainstreams, VI and MCMC. However, since larger dataset and more complex neural network are popular nowadays, the traditional approximate inference techniques are difficult to scale to large datasets and complex architectures. Therefore modern approximate inference techniques need to be adopted. In this section, two modern approximate inference techniques for BNN including dropout variational inference[GG16] and scalable Laplace approximation[RBB18] are introduced.

2.1.1 Introduction

Let \mathcal{D} denote an observable dataset consisting of a set of input and output pairs (\mathbf{x}_i, y_i) , that is $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\} = \{(\mathbf{x}_i, y_i)_{i=1}^N\}$, where $x_i \in \mathbb{R}^D$ and $y_i \in \mathcal{Y}$, \mathcal{Y} is the set of labels, and f^ω denotes one function parameterized by ω . Here in this case the model is a neural network with its parameters ω consisting of weights $\mathbf{W}_{1:L}$ and biases $\mathbf{b}_{1:L}$ for L layers.

In supervised learning, our goal is to learn a distribution of the output conditioned on the input, $p(\mathbf{y}|\mathbf{x})$ that can explain the underlying data distribution based on the observables \mathcal{D} which are samples drawn from the underlying data distribution. Then this 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:

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

where \mathbf{d} is possible class of random variable \mathbf{y} , and \mathbf{f}_d^w represents the probability that class \mathbf{d} is assigned to \mathbf{y} .

In the regression case, the likelihood is the Gaussian distribution:

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

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

As mentioned above, learning aims 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 the prior distribution of the model parameters $p(\omega)$. The probability distribution of the model parameters is updated from the prior into the posterior after observing the 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 the so-called model evidence or marginal likelihood, whose integration is always intractable when the dimensionality of parameters is too high.

After obtaining the posterior distribution over the model parameters, predictions can be obtained by marginalizing the likelihood of unseen input points such as \mathbf{x}^* over the model parameters, which leads to a predictive distribution over the output:

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

Explaining the difference between the deterministic neural network and the BNN can help understanding the mechanism of BNN better. In Fig.2.1, graphical model is used to express these two kinds of neural network, in which the solid point denotes the deterministic variable, the circle denotes the random variable, while the shaded circle denotes the observed random variable. The plate notation denotes \mathbf{N} observed data pairs (\mathbf{x}, \mathbf{y}) . Fig.2.1 demonstrates that parameters ω in the deterministic neural network is normal variable which is a point estimate obtained by MAP:

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

On the other hand, in BNN, the model parameters ω are random variables with the distribution function parameterized by θ . To note that, for explanation of concept here, no simplified assumptions are made on the function parameterized by θ , which means

it can be an arbitrary function. This distribution over model parameters can be inferred based on the Bayes' rule in Eq.2.3. However, for a model with a large number of parameters, it's hard to compute the model evidence (denominator of r.h.s of Eq.2.3) with integration tractably. An approximation needs to be employed to solve this problem. It's worth noting that if we choose the approximate distribution to be a delta function, $q_{\theta}(\omega) = \delta(\omega - \theta)$, the corresponding BNN is recovered into the deterministic neural network. Because the computation of model evidence $p(\mathbf{Y}|\mathbf{X})$ is always intractable in a complex model and a large dataset, approximate inference needs to be applied, which is briefly introduced in the next section.

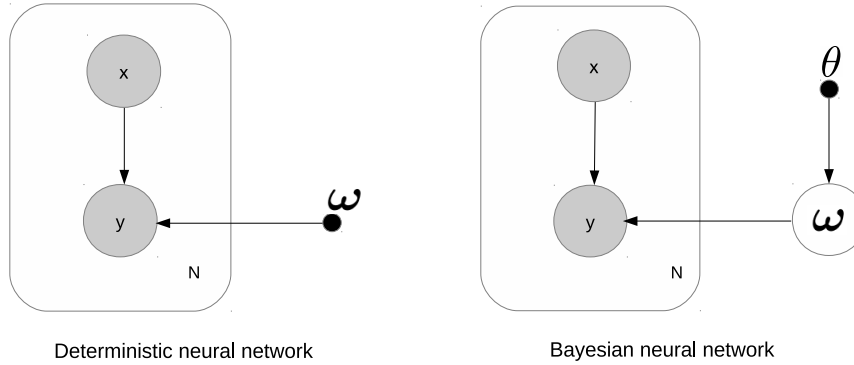


Figure 2.1: Differences between parameter estimation of the deterministic neural network and the BNN. The solid point denotes deterministic variable, and the circle denotes random variable, while the shaded circle denotes observed random variable. The plate notation denotes N observed data pairs (x, y) .

2.1.2 Dropout variational inference

VI casts inference into optimization by minimizing the KL-div between the approximate posterior distribution and the real posterior distribution. However, there is no analytical definition of this KL-div if the real posterior distribution is unknown. A lower bound can be derived with Jensen's inequality for tractable optimization, which is the so-called evidence lower bound (**ELBO**) which bounds the log marginal likelihood. And from that, the marginal likelihood is the sum of **ELBO** and KL-div between the approximate posterior

and the real posterior. The derivation is given in the following:

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

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

The $\log(p(Y)|X)$ can be obtained by calculating the sum of **ELBO** and KL-div between the approximate posterior $q_\theta(\omega)$ and the real posterior $p(\omega|X, Y)$:

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

When maximizing the **ELBO** w.r.t the parameter of approximate posterior θ , which is also called **variational parameter**, it is equivalent to minimizing the KL-div because the log marginal likelihood is not a function of θ . Then a well-defined objective **ELBO** comes out, 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 the regularization term which ensures the approximate posterior does not deviate too far from the prior. Now the inference problem has been cast into an optimization problem:

$$\theta^* = \underset{\theta}{\operatorname{argmin}}\{KL(q_\theta(\omega)||p(\omega|X, Y))\} \tag{2.8}$$

which is equivalent to

$$\theta^* = \underset{\theta}{\operatorname{argmax}}\{\mathbf{ELBO}\} \tag{2.9}$$

However, there are still difficulties solving this optimization problem with this objective. The first one is to deal with large size datasets, which induce too much computations in the

expected log likelihood term. [Gra11] shows that this can be solved by data sub-sampling which is similar as the stochastic optimization. Another one is the need of the derivatives of **ELBO** w.r.t. the variational parameters θ . Since the model parameters are samples from the approximate posterior, a good estimator for the derivatives is required which are introduced in the next subsection.

Dropout variational inference

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

The mechanism of dropout is straightforward, each units of specific layer is multiplied by a random variable under Bernoulli distribution with p as its parameter, where $1 - p$ is the so-called dropout rate. In each iteration during training, the 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 within the same iteration. In the testing, the sampling phase is turned off, only one forward propagation is needed to produce predictions. Normally, in order to avoid rescaling weights in the testing, which is used to keep the output magnitudes in the same scale when dropout is off, rescaling of the output of the dropout is always done in training. In Fig.2.2, there are two figures showing how the dropout is turned on and off.

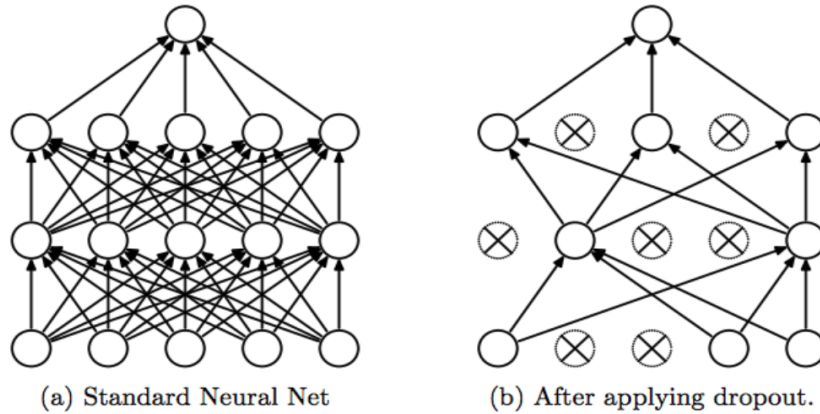


Figure 2.2: Schematic illustration of dropout in a three-layer fully connected neural network. (a) the standard neural network with all neurons switched on. (b) after applying dropout only a subset of neurons are used[SHK⁺14].

Bayesian interpretation of dropout: In VI, the goal is to minimize the KL-div between the approximate posterior and the real posterior over the model parameters. This objective is equivalent to the maximization of the 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 Fig.2.3, in which we define the hidden layer as the first layer and output layer as the second layer and assume a fully factorized Gaussian prior over parameters.

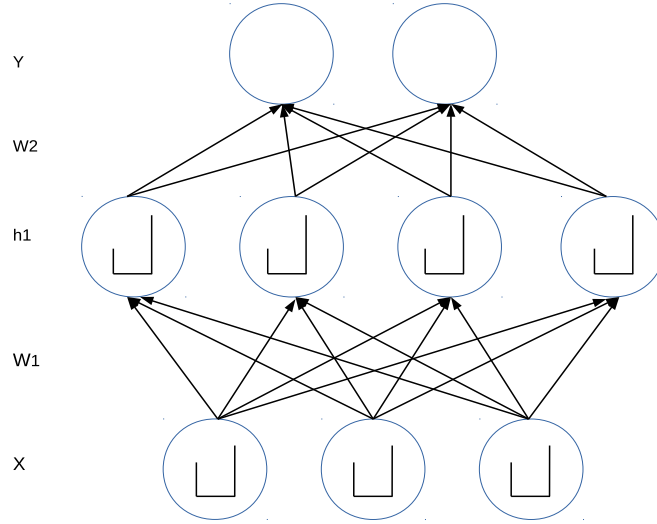


Figure 2.3: An example of two layer neural network with dropout, a Bernoulli random variable is imposed on each unit of input layer and hidden layer.

Approximate distribution: Fig.2.3 shows a simple example of DNN with dropout, in which the following variables are defined: $\mathbf{y} \in \mathbb{R}^{m \times D_2}$ as the output layer, $\mathbf{x} \in \mathbb{R}^{m \times D_0}$ as the input, $\mathbf{h}_1 \in \mathbb{R}^{D_1}$ as the response of the hidden layer, where m represents the number of data instances, and D_i is the dimensionality of the i -th layer, where the 0-th layer represents input layer and $i \in \{1, \dots, L\}$, $L = 2$. Further $\omega = \{(W_i)_{i=1}^L\}$ is defined as model parameters, and $\epsilon_i \in \mathbb{R}^{D_{i-1}}$ as the Bernoulli distributed random vector parameterized by $\mathbf{p}_i \in \mathbb{R}^{D_{i-1}}$, for i -th layer, which we call **keeping rate** in this report. The corresponding dropout rate is equal to $1 - \mathbf{p}_i$. It is assumed that a fully factorized Gaussian prior distribution is placed over the weights. These assumptions can easily be generalized and transferred to a multi-layer neural network easily.

In naive version of the dropout, each element of vector \mathbf{p}_i has the same value as \mathbf{p}_i , which means that there is only one keeping rate for each layer. In the following, if there is no further specification, \mathbf{p}_i and \mathbf{p}_i are used interchangeably depending on the context. Because the weights in weight matrix are treated as random variables, $\mathbf{M}_i \in \mathbb{R}^{D_{i-1} \times D_i}$ is used to denote position of non-zero element in Bernoulli distribution for \mathbf{W}_i . Additionally, the 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, also called homogeneous coordinate. Another assumption is that the approximate weight distribution is factorized over the layer, which yields

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

To start with the formulation of dropout, the likelihood of output conditioned on input with softmax scores of neural network is modeled in the classification case as:

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}, \omega) &= \sigma((\mathbf{h}_1 \odot \epsilon_2) \mathbf{M}_2) \\ &= \sigma(\mathbf{h}_1 (\text{diag}(\epsilon_2) \mathbf{M}_2)) \\ &= \sigma(\mathbf{h}_1 \mathbf{W}_2) \\ &= \sigma(\phi(\mathbf{x} (\text{diag}(\epsilon_1) \mathbf{M}_1) + \mathbf{b}_1) \mathbf{W}_2) \\ &= \sigma(\phi(\mathbf{x} \mathbf{W}_1) \mathbf{W}_2) \end{aligned} \tag{2.10}$$

where \odot is the Hadamard product (element-wise product), $\sigma(\mathbf{a}_j) = \frac{\exp(\mathbf{a}_j)}{\sum_k \exp(\mathbf{a}_k)}$ the softmax function, and $\phi(\cdot)$ the element-wise non-linear activation function such as the rectified unit function (Relu).

From the equation above, it can be derived:

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

which means that the weight matrix \mathbf{W}_i is a matrix of random variables 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 Eq.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

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

with

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

As can be seen, the approximate posterior over the model parameters puts the same Bernoulli distribution over the input dimension of the weight matrix, which is the dimension of row in this example. Meanwhile each element of the same row is multiplied with same realization of the random variable but with the different non-zero position, which is corresponding to the different expectation of each row element and thus induces the correlations between the row elements. To make this definition more clear, based on the computation of expectation and variance of the Bernoulli distribution, the first and second moment of the approximate distributed random variables can be written down like:

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

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

The covariance matrix of the parameter matrix is:

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

where

$$\begin{aligned} j &= l * D_{i-1} + q, \\ k &= m * D_{i-1} + n, \end{aligned}$$

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

From Eq.2.14, it can be seen 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 the covariance matrix of the parameter matrix \mathbf{W}_i w.r.t. the approximate posterior $q_{\theta}(\mathbf{W}_i)$, it is known that the samples of row vectors in the weight matrix are drawn independently and thus the covariance between rows are zero. On the other hand, the 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 the weights within the same row are induced. Therefore, weights within the same row can be learned by fitting this approximate distribution to the real posterior. This means that the approximate posterior

family has more flexibility to approximate the real posterior when compared with other common approximate posterior distribution family such as fully factorized Gaussian that assumes distribution for each parameter is independent.

Training objective: Up to now, only the approximate distribution of dropout inference is explained. As mentioned in introduction section of this chapter, it is known that by performing an optimization of **ELBO** w.r.t. variational parameters, a good approximation to the true posterior distribution over parameters can be obtained. This posterior can be used in testing to produce more reliable uncertainty estimation. In the following, training a DNN with dropout is equivalent to optimize the **ELBO** w.r.t. the variational parameters is illustrated.

At first, the training objective of DNN with dropout is defined as the cross entropy between the predictive distribution and the target distribution plus L2 regularization:

$$L_{dropout} = \sum_{i=1}^N [-\log(p(y_i|x_i, \omega))] + \lambda \left(\sum_{i=1}^L ||W_i||^2 \right) \quad (2.15)$$

where N represents the size of entire dataset, λ is L2 regularization coefficient. The likelihood over the entire dataset should be maximized w.r.t. the model parameter ω . This estimation of the model parameters is known as Maximum Likelihood Estimation (MLE). Equipped with L2 regularization and Gaussian prior over parameters, a MAP estimation is attained by minimizing Eq.2.15.

Normally gradient descent is used to tune the parameters during training, which requires first derivative of the objective w.r.t. the model parameters ω . Since large size of dataset is ubiquitous, which means N in Eq.2.15 is too large, exact gradients of whole batch (entire training set) is hard to calculate exactly and efficiently. Therefore many mini-batches of data are used to estimate gradients. This is the so called Stochastic Gradient Descent (SGD). Apart from making computation tractable, the noise of gradients estimations in each mini-batch is helpful for the optimization procedure to escape from the poor local minimum. The expression of gradients required in each iteration of SGD is given in the following:

$$\frac{\partial L_{dropout}}{\partial \omega} = \frac{N}{K} \sum_{i \in S} \left[-\frac{\partial \log(p(y_i|x_i, \omega))}{\partial \omega} \right] + \lambda \frac{\partial (\sum_{i=1}^L ||W_i||^2)}{\partial \omega} \quad (2.16)$$

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

On the other hand, when looking at the **ELBO** in Eq.2.6, to maximize **ELBO** is equivalent to minimize negative **ELBO** with SGD w.r.t. the variational parameter θ . The

gradients are computed with the following expression:

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

In order to calculate two terms in Eq.2.17, re-parameterization trick[KW13] needs to be introduced for the first term and the condition called KL condition[Gal] is introduced for the second term.

Re-parameterization trick: when taking a close look at the first term in Eq.2.17, it's required to compute the gradients of expectation w.r.t. the parameters of distribution to which this expectation is subject. That means, the gradients of variational parameters are random because our objective is evaluated with the realizations of these parameters. Fortunately, there are different approaches to estimate the gradients of variational parameters such as score function or REINFORCE estimator[Wil92], re-parameterization trick [KW13] and so on. As is stated in [KW13], the re-parameterization trick has lower variance than the score function estimator. There is a quick recap of this technique given and it can be seen that it is already a built-in part in DNN with dropout.

To identify the problem, a more general form is given in the following, the expression of derivatives of expectation w.r.t. the parameters upon which the distribution is defined:

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

Assuming that $p_\theta(x)$ can be re-parameterized as $p(\epsilon)$ which is a parameter-free distribution such that the 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 the estimator of the gradients w.r.t. the distribution parameters θ is as follows, the derivation of it can be referred in A.2:

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

From practical point of view, the expectation of last line in Eq.2.19 can be approximated with the Monte Carlo integration. In the dropout inference, it is known that if the dropout

rate is fixed to $\mathbf{1} - \mathbf{p}_i$ in Eq.2.12 during training, ϵ_i in Eq.2.12 is a parameter-free random variable and $\mathbf{g}(\cdot)$ is a differentiable function w.r.t. its input arguments \mathbf{M}_i . Now the variational parameters θ only contain $\{(\mathbf{M}_i)_{i=1}^L\}$, which are exactly the weights to be learned in training DNN with dropout. As a consequence, if the gradient of **ELBO** w.r.t. $\{(\mathbf{M}_i)_{i=1}^L\}$ can be computed, which is the first term in Eq.2.17, it is equivalent to calculate the gradients of the dropout loss w.r.t. the model parameters ω which is the first term in Eq.2.16.

KL condition: The second term in Eq.2.17 is proved to be equivalent to the second term in Eq.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

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

To establish this condition, we need to make a approximation of the approximate posterior distribution $q_\theta(\omega)$ in Eq.2.11, where we approximate $q_\theta(\mathbf{W}_i|\epsilon_i)$ in Eq.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 Gaussian distributions with small standard deviations, and one component with mean as zero:

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

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-div between mixture of Gaussian and a single Gaussian. Reader who has interest can refer to derivation of KL-div between mixture of Gaussian and single Gaussian in [Gal]. KL condition is given in the following:

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

with

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

for large enough D_i . If we **fix dropout rate** in training and compute the gradients of KL divergence w.r.t. variational parameters $\theta = \{(M_i)_{i=1}^L\}$, which is now exactly the same as the model parameters $\omega = \{(W_i)_{i=1}^L\}$. 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. Eq.2.16):

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

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

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

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

where $\omega \sim q_\theta(\omega)$ and $\hat{\omega}_i$ is one of K realizations of ω , which is equivalent to turning on dropout in testing time.

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 contain $\{(M_i)_{i=1}^L\}$, which are equivalent to $\omega = \{(M_i)_{i=1}^L\}$ in DNN with dropout. Therefore optimizing any DNNs with dropout is equivalent to perform approximate inference

on BNN with the same architecture. Thus the optimal weights learned during the optimization of DNN with dropout are the same as the optimal variational parameters of the corresponding BNN.

However, selecting dropout rate before training is not a trivial task for several reasons. Firstly, as shown in [SHK⁺14], different dropout rates have different influences on model capacity and thus model performance. To choose an optimal dropout rate manually requires repeating tedious experiments which induces a lot of time consumptions and computation efforts. Secondly, if we want the model not only to achieve satisfied accuracy but also to produce reliable uncertainty estimation. The dropout rate should matter a lot because it can influence the flexibility of approximate distribution family from the perspective of its 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 is not easy and intuitive to estimate gradients of expectation w.r.t. parameters of a discrete distribution. Before when the approximate distribution is continuous, gradients can be estimated 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 or categorical distribution, 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, which is called concrete distribution, categorical re-parameterization or Gumbel-softmax trick. In this trick, "max" operation in Gumbel-Max trick is replaced by "softmax" function, which can yield 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 for Bernoulli distribution: Gumbel-Max trick [MTM14] is shown in Fig.2.1.2, 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 label representation, that is $d \in \{0, 1\}^n$ and $\sum_{i=1}^n d_i = 1$. The Gumbel-Max trick proceeds as follows(cf. Fig.2.1.2):

- 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.

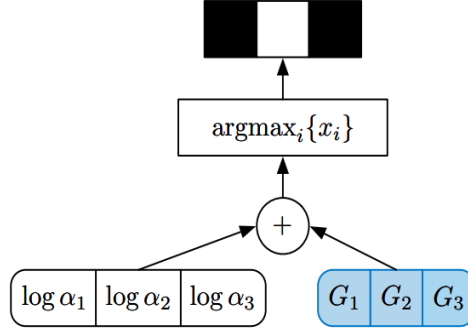


Figure 2.4: Example of Gumbel-max trick for drawing samples from a discrete distribution whose random variable has 3 states and $\{\alpha_i\}_{i=1,2,3}$ as class parameters representing the possibility of occurrence of that class. $\{G_i\}_{i=1,2,3}$ are i.i.d Gumbel($\mathbf{0}, \mathbf{1}$) [MMT16].

As observed in Gumbel-Max trick, the sampling step is re-parameterized 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. This function is differential w.r.t. distribution parameters. With this replacement, a continuous approximation to one-hot encoding vector \mathbf{d} on the simplex $\Delta^{n-1} = \{\mathbf{x} \in \mathbb{R}^n | \mathbf{x}_k \in [0, 1], \sum_{i=1}^n \mathbf{x}_i = 1\}$ is obtained:

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

The sampling steps are similar to Gumbel-Max trick, but with smoothed **softmax** operation instead of **argmax** (cf. Fig.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.

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:

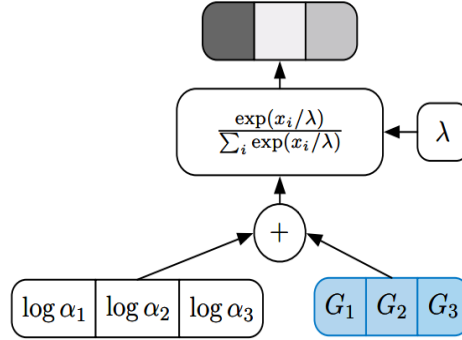


Figure 2.5: Example of continuous approximation to Gumbel-max trick for drawing samples from a discrete distribution whose random variable has 3 states and $\{\alpha_i\}_{i=1,2,3}$ as class parameters representing the possibility of occurrence of that class. $\{G_i\}_{i=1,2,3}$ are i.i.d Gumbel($\mathbf{0},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

$$\begin{aligned}\mathbb{P}(\mathbf{d}_1 = 1) &= \mathbb{P}(\log \alpha_1 + \mathbf{G}_1 > \log \alpha_2 + \mathbf{G}_2) \\ &= \mathbb{P}(\log \alpha_1 - \log \alpha_2 + \mathbf{G}_1 - \mathbf{G}_2 > 0)\end{aligned}$$

where $\mathbf{G}_1 - \mathbf{G}_2 = \log(\text{Uniform}(\mathbf{0}, \mathbf{1})) - \log(1 - \text{Uniform}(\mathbf{0}, \mathbf{1}))$, which is the inverse cumulative density function of Logistic(0,1). Then we can infer the re-parameterization of continuous approximation of sample $\mathbf{x} \in [\mathbf{0}, \mathbf{1}]$ from Bernoulli distributed random variable with \mathbf{p} as parameter as follows:

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

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

In order to make explanation more clear and validate the continuous approximate re-parameterization, one plot (Fig.2.6) is given to show 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 increasing step of 0.1. It can be seen that the empirical average of samples drawn from Eq.2.25 is close to expectation of Bernoulli distribution of corresponding parameter. Additionally, one single sample in the figure is shown. As can be seen, most of samples are lying on the boundary of range $[\mathbf{0}, \mathbf{1}]$ and only few of them lie in the interior.

Dropout regularization: Because keeping rate of dropout is being optimized in concrete dropout, the variational parameters $\boldsymbol{\theta}$ contain both $\{(M_i)_{i=1}^L\}$ and $\{(\mathbf{p}_i)_{i=1}^L\}$. In order to compute gradients of ELBO in Eq.2.17, categorical re-parameterization is employed to estimate the gradients w.r.t. Bernoulli distribution parameter \mathbf{p}_i in the first

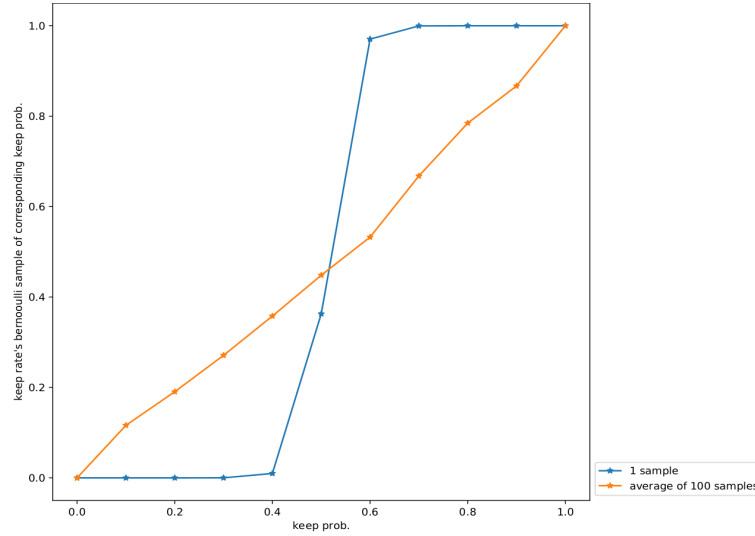


Figure 2.6: One sample and average value of 100 samples drawn from continuous approximation of Bernoulli distribution with parameter $\mathbf{p} = [0.1, 0.2, \dots, 1.0]$ and temperature $\lambda = 0.1$.

term, for the second term, we could not ignore the term with \mathbf{p}_i in KL condition Eq.2.21, which is the entropy term. Consequently, unlike Eq.2.22, the KL-div term should be:

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

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, it can be seen 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 regularization 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.

2.1.3 Laplace approximation

Introduction

In this section, another method to approximate the real posterior distribution over the weights is introduced, which is called Laplace approximation[Bis06]. It requires only the MAP point estimation of the model parameters, thus there is no need to re-train the existing model which has already obtained an point estimate. This characteristic makes this approach more appealing because it is not required to modify the model and thus can be applied to any existing models. The idea behind Laplace approximation is to put a Gaussian approximation on the MAP point estimate of the model parameters, which is one mode of the true posterior. The reasons why we choose this approach are as following:

- it has the easy compatibility to the existing network. To perform Laplace approximation, we only need one point estimation of the model parameters which is already available for all DNNs. There is no need to modify the training phase which saves a lot of computation and time.
- it can capture the relationship between model parameters. As mentioned in the first chapter, most of approximation approaches make an assumption that parameters are independent to each other for reasons like simplicity, scalability as well as efficient computation. This could be a quite strong restrictions on the approximation which would lead to bad performance.

In the following, the basic idea of Laplace approximation is introduced and further its scalable version for DNN based on Kronecker factor approximation is explained.

Assume there is already a point estimate of the model parameters via MAP or MLE:

$$\begin{aligned}
 \omega^* &= \underset{\omega}{\operatorname{argmax}} \{p(\omega|Y, X)\} \\
 &= \underset{\omega}{\operatorname{argmax}} \left\{ \frac{p(Y|X, \omega)p(\omega)}{p(Y|X)} \right\} \\
 &= \underset{\omega}{\operatorname{argmax}} \{p(Y|X, \omega)p(\omega)\}
 \end{aligned} \tag{2.27}$$

After taking a second order Taylor expansion of the logarithm of of the posterior distribution $p(\omega|Y, X)$ around its mode ω^* (cf. A.3), by assuming that the prior of weights is uniform, the negative Hessian can be expressed as follows:

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

When exponentiating the second order Taylor expansion of log-posterior, the following

Gaussian-like functional form can be attained:

$$p(\omega|Y, X) \propto p(\omega^*|Y, X) \exp\left\{-\frac{1}{2}(\omega - \omega^*)^T H(\omega - \omega^*)\right\} \quad (2.28)$$

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

If Gaussian prior is chosen, it's easy to know that the second order derivative of the prior distribution term $\frac{\partial^2 p(\omega)}{\partial \omega^2}$ 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, the negative Hessian of log likelihood is defined as \hat{H} as follows:

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

However, for a large training set, it is always infeasible to analyze the gradients or Hessian exactly. To resolve this, normally the expectation of gradients or Hessian is estimated with a mini-batch of data. It means that the Hessian can also be estimated with the empirical average Hessian computed in mini-batches:

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

where N is size of training samples, and

$$\mathbb{E}_{p(Y, X)}[\hat{H}] \approx -\frac{1}{K} \sum_k \left[\frac{1}{M} \sum_i \frac{\partial^2 \log(p(y_{ik}|x_{ik}, \omega))}{\partial \omega^2} \right] \quad (2.29)$$

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

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

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

Therefore the Fisher information matrix can be used as replacement of the expected Hessian for the log likelihood term. The derivation of equivalence between expected Hessian and Fisher matrix is straightforward. We define Fisher matrix F in the following:

$$\begin{aligned} \hat{F} &= \mathbb{E}_{p(Y, X)} \left[\frac{\partial}{\partial \omega} \log p(Y|X, \omega) \frac{\partial}{\partial \omega} \log p(Y|X, \omega)^T \right] \\ &\approx \frac{1}{K} \sum_k \left[\frac{1}{M} \sum_i \left(\frac{\partial}{\partial \omega} \log p(y_{ik}|x_{ik}, \omega) \frac{\partial}{\partial \omega} \log p(y_{ik}|x_{ik}, \omega)^T \right) \right] \end{aligned} \quad (2.30)$$

We include the derivation of equivalence between negative expected Hessian and Fisher matrix for case that parameter is scalar in appendix A.1, which can be extended to case for vector easily. The relationship between Fisher information matrix and expected negative Hessian when considering the entire training set is as follows:

$$\hat{H} \approx N \hat{F} \quad (2.31)$$

and after taking the Gaussian prior into account:

$$H \approx N \hat{F} + \tau I \quad (2.32)$$

To compute the empirical average of outer product, large storage space should be used to save values with amount quadratic to the number of model parameters. More than that, the computation complexity of matrix inversion is cubic of dimensionality of the matrix which is the number of parameters of the whole network. This is also infeasible because the number of parameters in DNN can exceed million easily nowadays. Therefore an efficient approximation for Fisher matrix is required.

Scalable Laplace approximation for neural network

In order to mitigate the computational burden in Laplace approximation above, Kronecker-factored approximation curve(KFAC) in [MG15] can be used to approximate the Fisher information matrix \mathbf{F} , with which Laplace approximation can be performed for neural network[RBB18]. 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 Fisher information matrix is only required after training, this little expensive computation is need to perform only once.

With the notations used in previous subsection, where dropout variational inference is introduced, the derivation of the approximation is given in the following. For clarity, we repeat the definitions of existing notations and new defined ones here. A DNN transforms its input $\mathbf{a}_0 = \mathbf{x}$ to an output $\mathbf{a}_L = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) = \boldsymbol{\sigma}(\mathbf{x}, \boldsymbol{\omega})$ 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:

$$\mathbf{s}_i = \mathbf{a}_{i-1} \mathbf{W}_i \text{ and } \mathbf{a}_i = \phi(\mathbf{s}_i) \quad (2.33)$$

We define $\tilde{\boldsymbol{\omega}}$ to be the vector consisting of all of the network parameters concatenated together, i.e. $\tilde{\boldsymbol{\omega}} = [\text{vec}(\mathbf{W}_1)^T \text{vec}(\mathbf{W}_2)^T \dots \text{vec}(\mathbf{W}_L)^T]^T$, where vec is the operator

which vectorizes matrices by stacking their columns together. For simplicity the following notations are defined:

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

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

$$\partial \tilde{\mathbf{w}} = [\text{vec}(\partial \mathbf{W}_1)^T \text{vec}(\partial \mathbf{W}_2)^T \dots \text{vec}(\partial \mathbf{W}_L)^T]^T \quad (2.34)$$

From Eq.2.30, $\hat{\mathbf{F}} = \mathbb{E}[\partial \hat{\mathbf{w}} \partial \hat{\mathbf{w}}^T]$ 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 $\text{vec}(uv^T) = v \otimes u$, 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 by rewriting it into:

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

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

$$\begin{aligned} \tilde{\mathbf{F}}_{i,j} &\approx \hat{\mathbf{F}}_{i,j} \\ &= \mathbb{E}[\mathbf{g}_i \mathbf{g}_j^T] \otimes \mathbb{E}[\mathbf{a}_{i-1} \mathbf{a}_{j-1}^T] \\ &= \mathbb{E}[\mathbf{G}_{i,j}] \otimes \mathbb{E}[\mathbf{A}_{i,j}] \end{aligned} \quad (2.35)$$

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 large diagonal block matrix, which means $\tilde{\mathbf{F}}_{i,j}$ is non-zero only for $i = j$. For i -th layer, the Fisher information matrix can be computed with Eq.2.35. 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 they are treated as covariance matrix of resulting Gaussian. That means each Gaussian has a Kronecker factored covariance matrix, 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:

$$\mathbf{W}_i \sim \mathcal{MN}(\mathbf{W}_i^*, \mathbb{E}[\mathbf{A}_i]^{-1}, \mathbb{E}[\mathbf{G}_i]^{-1}) \quad (2.36)$$

If Gaussian prior and scale of Fisher information matrix(cf. Eq.2.32) are taken into account, then the resulting posterior can be written in the following form:

$$\mathbf{W}_i \sim \mathcal{MN}(\mathbf{W}_i^*, (\sqrt{N}\mathbb{E}[\mathbf{A}_i] + \sqrt{\tau}\mathbf{I})^{-1}, (\sqrt{N}\mathbb{E}[\mathbf{G}_i] + \sqrt{\tau}\mathbf{I})^{-1}) \quad (2.37)$$

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.

2.2 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. 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.

2.2.1 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 $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and their labels $\mathbf{y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$. 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, the probability distribution of \mathbf{y} conditioned \mathbf{x} with model parameter θ can be defined as

$$p(\mathbf{y}|\mathbf{x}; \theta) = \frac{1}{Z(\mathbf{x}, \theta)} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{y}_c, \mathbf{x}_c, \theta) \quad (2.38)$$

where \mathcal{C} is the set of cliques of graph \mathcal{G} , and $Z(\cdot)$ is normalizer to make sure $\sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y}|\mathbf{x}; \theta) = 1$, which is also called partition function. θ 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, only factors with number of random variables up to two are considered in common 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 $\psi_i(\mathbf{y}_i, \mathbf{x}_i, \theta)$. 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 $\psi_{i,j}(\mathbf{y}_i, \mathbf{y}_j, \mathbf{x}_i, \mathbf{x}_j, \theta)$ characterized by edges in the graph. This kind of potential states the compatibility of two random variables $(\mathbf{y}_i, \mathbf{y}_j)$ being assigned to a certain pair of categories.

To make the explanation more clear, we illustrate a simple example in Fig.2.7, we ignore θ here and explain it later and use $\psi_{i,j}$ for $\psi_{i,j}(\mathbf{y}_i, \mathbf{y}_j, \mathbf{x}_i, \mathbf{x}_j)$ to make expression of conditional probability uncluttered:

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

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

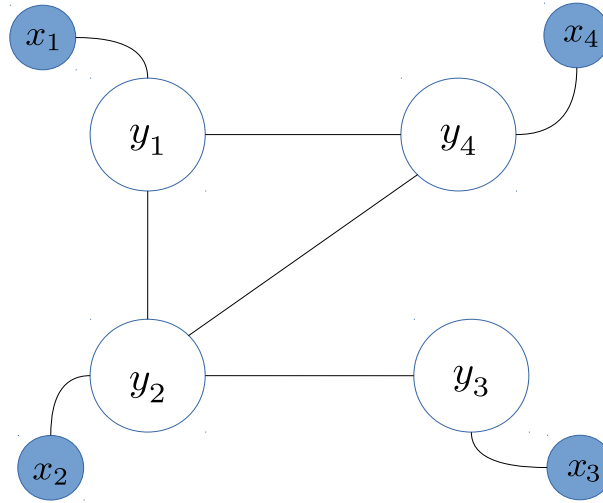


Figure 2.7: Simple example of conditional random field.

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:

$$\begin{aligned}
p(y|x; \theta) &= \frac{1}{Z(x, \theta)} \prod_{c \in \mathcal{C}} \exp(\langle \phi(x_c, y_c), \theta \rangle) \\
&= \frac{1}{Z(x, \theta)} \exp\left(\sum_{i \in V} \langle \phi_u(x_i, y_i), \theta_u \rangle + \sum_{(i,j) \in \mathcal{E}} \langle \phi_p(x_i, x_j, y_i, y_j), \theta_p \rangle\right)
\end{aligned} \tag{2.39}$$

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 only 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.2.39. This means that the dimension of the parameter in 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 also 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.

Regarding **unary feature**, we utilize the probability vector $\mathbf{p}(\mathbf{y}_i|\mathbf{x}_i)$ which is predictive likelihood from a discriminative classifier which is a deep neural network in this work.

Pairwise feature in this work is quite simple context cue in the scene which is defined as binary co-occurrence matrix $\mathbf{H}(\mathbf{y}_i, \mathbf{y}_j)$. If one category occurs with another category in the same scene, then the corresponding entry in the matrix is filled with 1, otherwise 0.

Although obtaining this pairwise feature requires knowing the ground truth labels of objects in the scene, it still makes sense in the application of this work. Because we use T-LESS dataset which will be introduced in the experiment part, each of scene images in this dataset contains multiple industrial-relevant texture-less objects, some of them are combination of objects of other categories. We want to simulate the situation that, we treat each scene as one product which is composed of multiple sub-parts. Because we know the product, and thus we also know the real labels of components of this product which is corresponding to the co-occurrence of objects in this scene and our goal is to classify the components with help of this contextual information.

Based on these two kinds of features, we have the formulation of likelihood function of CRF model for one scene as follows:

$$p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\mathbf{x}, \boldsymbol{\theta})} \exp(\theta_u \sum_{i \in V} p(\mathbf{y}_i|\mathbf{x}_i) + \theta_p \sum_{(i,j) \in \mathcal{E}} H(\mathbf{y}_i, \mathbf{y}_j)) \quad (2.40)$$

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

2.2.2 Learning

The most popular approach is estimate $\boldsymbol{\theta}$ is MLE. The likelihood for one scene is defined in Eq.2.40. Normally we have large size of scenes in training set where each of them contains multiple objects. We denote training set by $\mathcal{D} = [\mathbf{d}_1, \dots, \mathbf{d}_m]$, where each scene is represented by $\mathbf{d}_i = (\mathbf{x}^i, \mathbf{y}^i)$. 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 $\mathcal{G}_i = (\mathbf{V}_i, \mathcal{E}_i)$. 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:

$$\mathcal{NLL}(\theta; \mathcal{D}) = - \sum_{i=1}^m (\theta_u \sum_{j \in V_i} p(y_j^i | x_j^i) + \theta_p \sum_{(j,k) \in \mathcal{E}_i} H(y_j^i, y_k^i) - Z_l(x^i; \theta)) \quad (2.41)$$

where

$$Z_l = \log \left(\sum_{\mathbf{y}^i \in \mathcal{Y}} (\theta_u \sum_{j \in V_i} p(y_j^i | x_j^i) + \theta_p \sum_{(j,k) \in \mathcal{E}_i} H(y_j^i, y_k^i)) \right) \quad (2.42)$$

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

$$\theta^* = \underset{\theta}{\operatorname{argmin}} [\mathcal{NLL}(\theta; \mathcal{D})] \quad (2.43)$$

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 SGD, which requires to calculate the gradient of object w.r.t. the model parameters θ . The gradients can be computed as follows, where we use θ to represent parameters in order to make the expression clear:

$$\frac{\partial}{\partial \theta} \mathcal{NLL}(\theta; \mathcal{D}) = - \sum_{i=1}^m (\phi(x^i, y^i) - \mathbb{E}_{p(y|x^i; \theta)} [\phi(x^i, y)]) \quad (2.44)$$

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

$$\mathbb{E}_{p(y|x^i; \theta)} [\phi(x^i, y)] = \sum_{y \in \mathcal{Y}} p(y|x^i; \theta) \phi(y, x^i) \quad (2.45)$$

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 $p(y|x; \theta)$. 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.2.44 and Eq.2.45, to calculate the gradients in each step of SGD, we need to compute $p(y|x; \theta)$, 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 $\mathbf{Z}_l(\cdot)$, which is required in Eq.2.45 and to the marginal probability distribution of random variables in \mathbf{y} to be plugged in Eq.2.44 in the graph with cycles and exact solution in graph without cycle.

2.2.3 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 \mathbf{k} -dimensional vector $\mathbf{b}(\mathbf{y}_i)$ if the node is discrete random variable where \mathbf{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 its neighbors. A message $\mathbf{m}_{ij}(\mathbf{y}_j)$ sent from node \mathbf{y}_i to node \mathbf{y}_j encodes belief of \mathbf{y}_i about what class node \mathbf{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:

$$\mathbf{m}_{ij}(\mathbf{y}_j) = \sum_{\mathbf{y}_i \in \mathcal{Y}_i} \psi_{ij}(\mathbf{y}_i, \mathbf{y}_j, \mathbf{x}_i, \mathbf{x}_j, \theta_p) \psi_i(\mathbf{y}_i, \mathbf{x}_i, \theta_u) \prod_{\mathbf{y}_k \in \mathcal{N}(\mathbf{y}_i) \setminus \mathbf{y}_j} \mathbf{m}_{ki}(\mathbf{y}_i) \quad (2.46)$$

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

$$\mathbf{b}(\mathbf{y}_i) = \frac{1}{\mathbf{Z}_i} \psi_i(\mathbf{y}_i, \mathbf{x}_i, \theta_u) \prod_{\mathbf{y}_j \in \mathcal{N}(\mathbf{y}_i)} \mathbf{m}_{ji}(\mathbf{y}_i) \quad (2.47)$$

where \mathbf{Z}_i is the normalizer of unnormalized belief of node \mathbf{y}_i to ensure the sum of marginal probability of \mathbf{y}_i is 1.

Chapter 3

Technical Approach

In this chapter, we will explain the approaches employed in this work. Firstly, confronting with the limitations of concrete dropout, we propose another version of concrete dropout, which is called multiple dropout(multi-drop). Secondly, we introduce the modified architecture of ResNet50 in order to incorporate different inference techniques for BNN with this powerful feature extraction network. Thirdly, combination of BNN and CRF in this work is introduced, where CRF works as a post-processor to capture the relationship between random variables. Last but not least, we explain the approach combining aforementioned techniques for continuous learning in the last section.

3.1 Multiple dropout

From Fig.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. Fig.3.3), which means random vector $\mathbf{p}_i = [p_i^k]_{k=1}^{D_i-1}$. The difference between multiple dropout and concrete dropout can be seen easily by comparing 3.3 and 2.3. While the first

term in gradients computation stays the same, the second term should be modified to:

$$\begin{aligned} \frac{\partial KL(q_\theta(\omega)||p(\omega))}{\partial \theta} &\approx \frac{\partial \sum_{i=1}^L \lambda ||M_i||^2 - \beta \mathcal{H}(p_i)}{\partial \theta} \\ &= \frac{\partial}{\partial \theta} \left(\sum_{i=1}^L \lambda ||M_i||^2 - \beta \sum_{k=1}^{D_{i-1}} (-p_i^k \log p_i^k - (1 - p_i^k) \log(1 - p_i^k)) \right) \end{aligned} \quad (3.1)$$

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. Fig.3.1), 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.

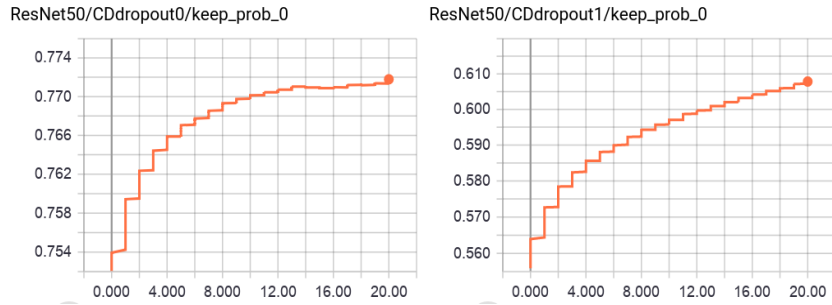


Figure 3.1: Changes of keep probability of the first two concrete dropout layers during training.

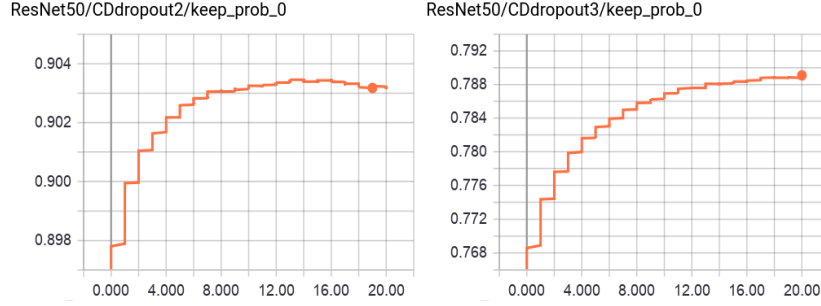


Figure 3.2: Changes of keep probability of the last two concrete dropout layers during training.

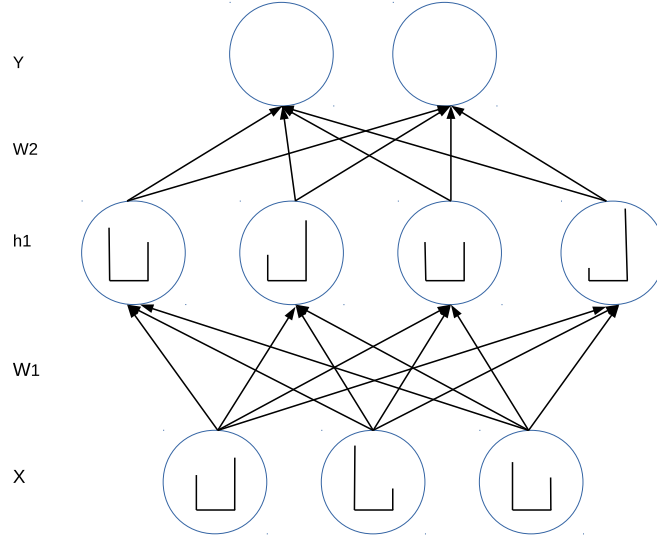


Figure 3.3: Different dropout rates for different hidden units in multi-drop.

3.2 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:

- Inserting dropout in layers initialized with pre-trained weights will destroy pre-trained features. Because we initialize major part of network with pre-trained weights, hence it would lead to significant drop of performance after fine-tuning if we insert dropout into them.
- 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.

Fig.3.4 shows the sketch of our modified network architecture. We can interpret major part of network, which do not have dropout inserted, as a deterministic feature extractor. Following this, there four layers with dropout inserted work as a probabilistic classifier. These four layers include one flatten layer with 2048 hidden units, and three fully connected layer with 1000 hidden units each. At the end, an output layer with size of number of categories is appended. The version of dropout inserted can be normal dropout, concrete dropout or multiple dropout. In this report, we only show results of concrete dropout and multiple dropout because normal dropout underperforms when compared with them.

In training, these two parts are trained together in order to achieve a better balance between uncertainty estimation and accuracy. As known to all, freezing feature extraction part can lead to worse performance, but it can also reduce the influence on predictive uncertainty from aleatoric uncertainty. To investigate this effect, we have done a ablation study which is shown in experiment part. In testing, we need to marginalize possible parameters according to posterior distribution. Layers with dropout should be turned on, sampled and run multiple times to perform Monte Carlo integration to approximate predictive distribution. The reduction of run time in testing can be achieved with different techniques, such as parallel computing or distillation. But this is out of the scope of this work.

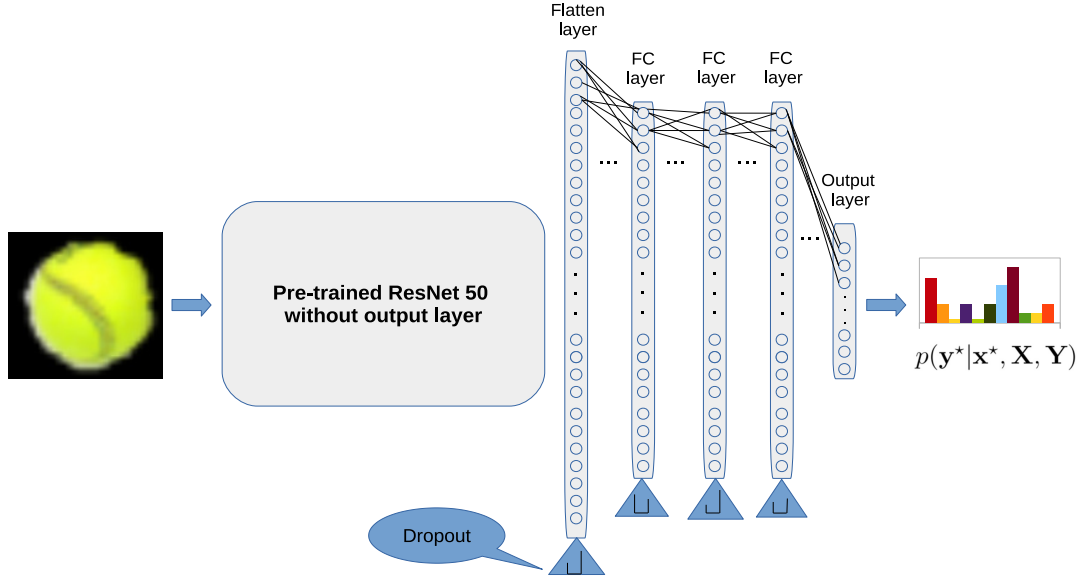


Figure 3.4: Modified network architecture of ResNet50.

3.3 Combination with CRF

After employing Bayesian neural network, the model can provide more reliable uncertainty estimation in a probabilistic manner. In other words, the output distribution can bring more information than before. It would be better to make use of this information with another probabilistic model. In this work, we explore to combine output from BNN and relationship between objects in scene via CRF. In detail, we use output distribution from BNN for unary feature and binary co-occurrence matrix for pairwise feature (introduced in 2.2.1). Fig.3.5 demonstrate the flow chart how BNN and CRF is combined together. As we can see, objects in one scene are firstly fed into BNN and their output distributions are fed into CRF as unary feature for each node. In the scene, their contextual relationship is represented by the binary co-occurrence matrix which is fed into CRF as pairwise feature. Then loopy belief propagation is used to infer the marginal distribution of each node.

One key factor in this case is that the pairwise feature is required to be designed by hand. The choice and design of it have significant impact on the improvement achieved by CRF. The information brought by pairwise feature should be complementary instead of contradictory to that of unary feature. Once this point is fulfilled, the more information unary feature has, the more improvement CRF can achieve.

One thing worth noting, in experiment part, we test this approach only on T-LESS dataset. Because to construct scene for objects in WRGBD dataset is non-trivial. The reason for this is that most objects there have not only similar appearances but also similar contexts.

Therefore it's hard to disambiguate miss-classification by incorporating their contextual information.

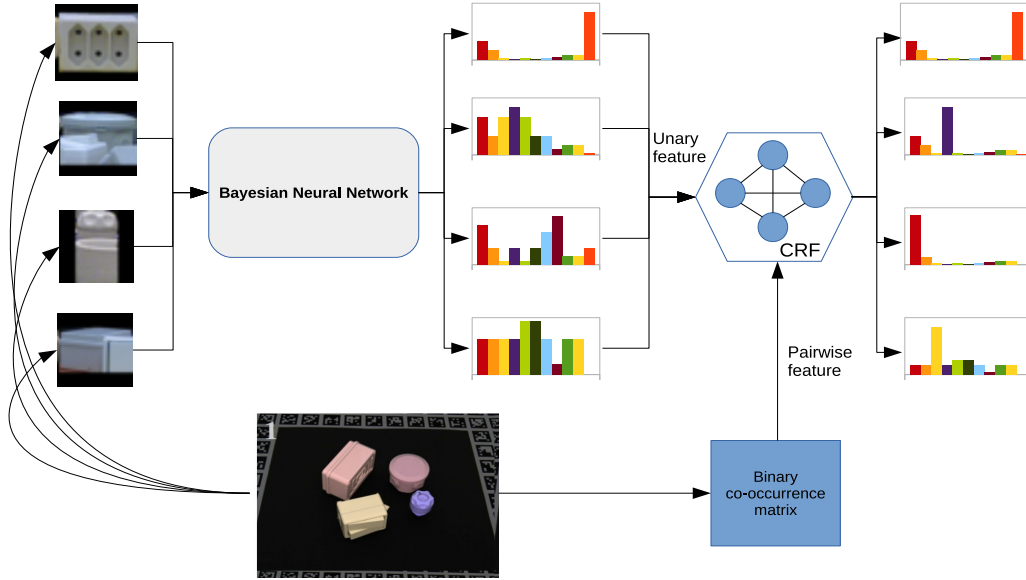


Figure 3.5: Combination of BNN and CRF.

3.4 Approach for continuous learning

With reliable uncertainty estimation, one of our goals is to enable the classifier or robot to learn continuously. The classifier should be introspective and express the confidence about its predictions. It's common that the test data in real environment does not have exactly same distribution as the training set, which induces a significant performance drop. In this situation, we expect our classifier to be able to adapt itself in real environment by fine-tuning itself with as little manual effort as possible.

This idea is visualized in Fig.3.6. The initialization phase includes training model with easily obtainable dataset, which can be public or synthetic dataset. Before deploying in the real environment, there is an adaptation phase which makes use of the aforementioned techniques. In this phase, the classifier should be able to adapt itself to the objects in real environment by fine-tuning with dataset collected with as little manual effort as possible. If the relationship between objects in real environment is complementary to the discriminative classifier and can be encoded well, CRF can be employed to capture the relationship to improve the performance further (cf. section 3.3). After that, the classifier is deployed to test data in real environment. In the experiment part, we evaluate this approach on both UniHB(WRGBD) and T-LESS(synthetic T-LESS).

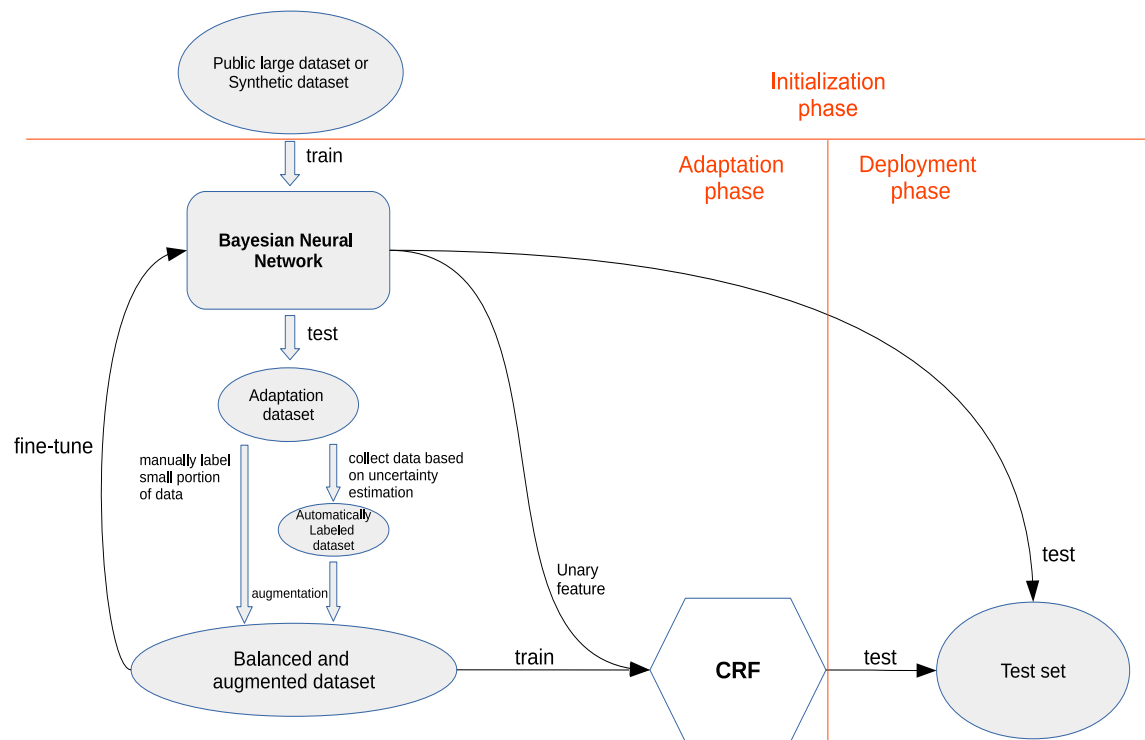


Figure 3.6: Approach for continuous learning.

Chapter 4

Experiment and Discussion

4.1 Background

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 some details of experiments such as the dataset, evaluation metrics. Models of deep neural nets are 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.

Regarding implementation of conditional random field, we employ C++ library from [RSGGJ15]. However, the weights there 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 modify the type of weights in this library in order to incorporate different types of weight defined in Eq. 2.40, which serve as coefficients of different potential functions instead of a linear mapping inside the potential functions. To note that the messages are initialized uniformly during loopy belief propagation and parameter are initialized with Gaussian prior with small standard deviation (0.01) around 0.

Besides, we perform optimization with stochastic gradient descent(SGD) with momentum and set size of mini-batch as 16, learning rate as $1e^{-4}$ and maximum number of iteration as 15×10^3 . Loopy belief propagation(LBP) is used for inference both during training and testing. The maximum number of iteration for computing messages in LBP is 100, and the convergence threshold is $1e^{-4}$.

4.1.1 Dataset

WRGBD[LBRF11]

This dataset is a public large-scale dataset of 300 household objects captured from multi-viewpoint. The objects are organized into 51 categories and 300 instance classes in hierarchical structure. There are multiple instances in different categories. 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. Fig.4.1). The total size of entire dataset is around $\sim 160.9K$. 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 in features space 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.

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. In this work, we call it UniHB dataset for simplicity. The size of data with elevation 45° is around $\sim 8.6K$ and that of data with elevation 30° and 60° is around $\sim 17.1K$.

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. Fig.4.2) that do not belong to the 51 class 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 classes 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 $\sim 6.0K$.

Original T-LESS[HHO⁺17]

This dataset features 30 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 captured systematically sampled from a view sphere, resulting in around $\sim 38.0\text{K}$ training images (cf. on the left of each thumbnail in Fig.4.3) and $\sim 10.0\text{K}$ test scene images (cf. on the right of each thumbnail in Fig.4.3). The training images depict objects in isolation with a black background, while the test images are from 20 table-top scenes with arbitrarily arranged objects placed on table. In this work we only consider the cropped images of objects in testing scenes, resulting a testing set with size $\sim 69.5\text{K}$. Nevertheless, the context information in the scene can be exploited with conditional random field which is showed in the experiment part.

Synthetic T-LESS

Because collecting data of real objects requires too high amount of manual efforts, our idea is to first train classifier on synthetic dataset which are easily obtainable. Automatically labeled data can be collected based on improved uncertainty estimation, which can then be used to train a more accurate classifier. For this end, we create a synthetic T-LESS dataset which contains multi-view images of objects rendered from 3D reconstructed CAD models with similar augmentations(cf. on the oright of each thumbnail in Fig.4.4). There is significantly large domain gap between the synthetic objects and real objects (cf. on the left of each thumbnail in Fig.4.4) such as texture and geometries on the surface. This domain gap can induce large performance drop when we test the model trained with synthetic objects on the test set with real objects.

Augmentation is employed for both original and synthetic T-LESS dataset. Because objects in test scenes are mostly occluded and have different backgrounds and lightness, this practical technique can improve the generalization of model and thus performance. In augmentation, the real and synthetic single objects are put onto different background images from VOC2012 dataset[EVGW⁺](cf. Fig.4.4) and augmented in different ways such as Gaussian noise, color changes, rotations and so on. One characteristic of the augmentation is that each background image can be treated as one scene, in which random number of objects are placed randomly, inducing large variety of occlusions and combinations of different objects. This tries to simulate the situation in test set, where different objects are placed randomly in each scene. More than that, the relationship between objects in one scene can be utilized by CRF for further performance gain.

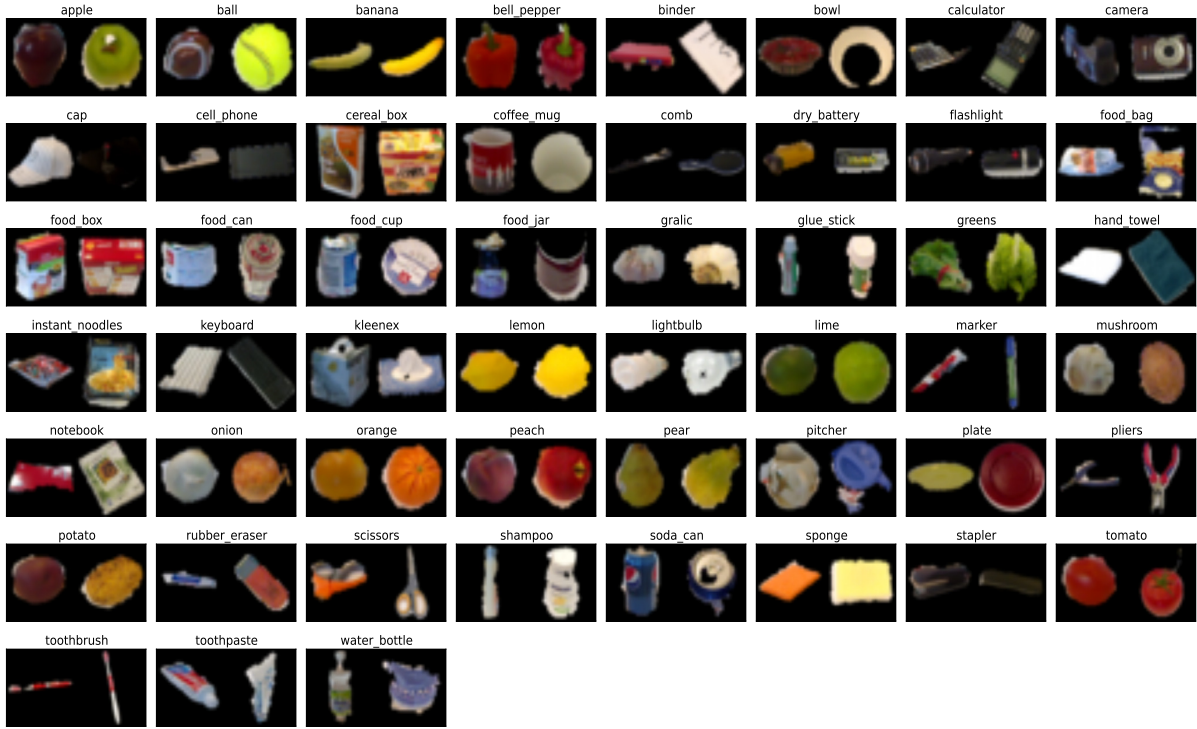


Figure 4.1: Example of masked images of objects from 51 categories in WRGBD and UniHB dataset. In each category, the left is from WRGBD and the right is from UniHB. We randomly pick one instance for the objects of WRGBD. We can see some light and appearance difference between objects in these two datasets.

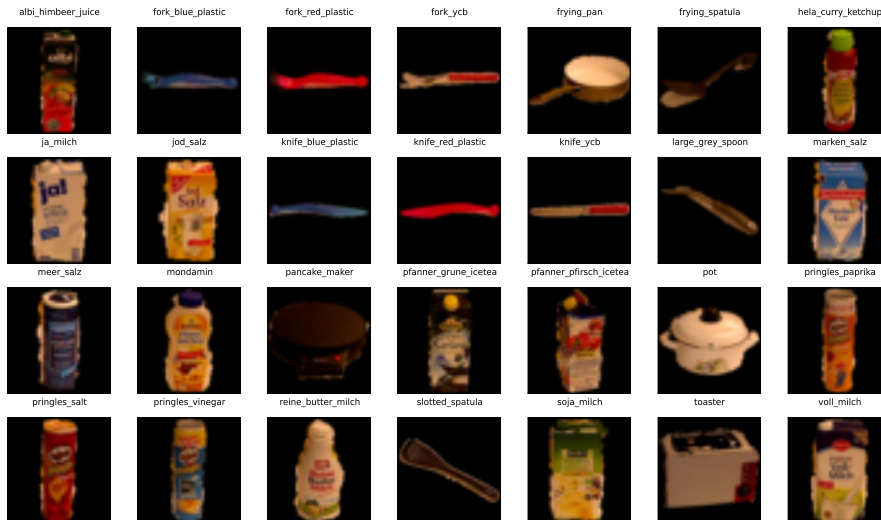


Figure 4.2: Example of masked images of objects from 28 categories which are not belonging to WRGBD categories, which are treated as OOD data samples.

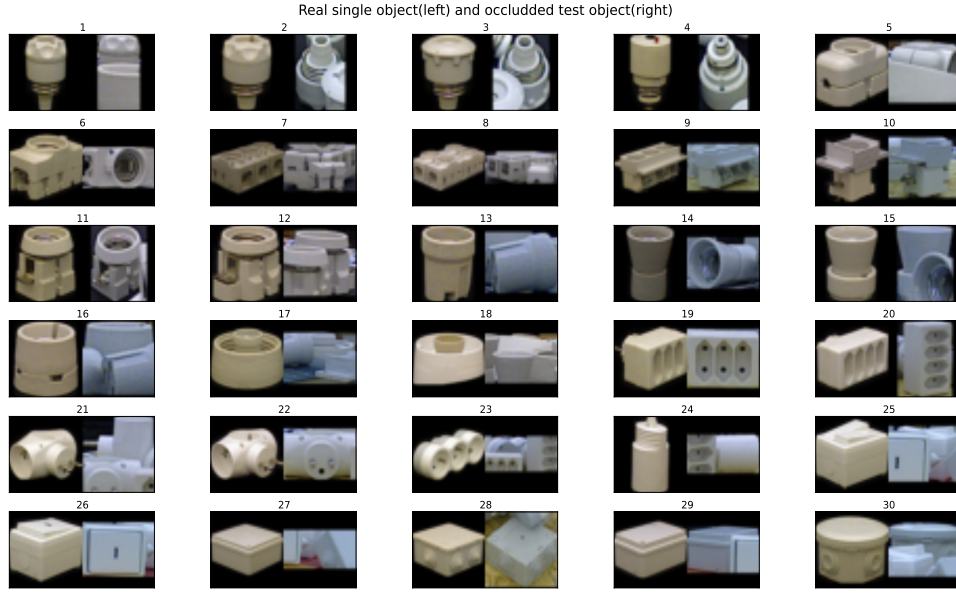


Figure 4.3: Example of object of each category in original training set and test set. In each thumbnail, the left image is real single object with black background and the right image is cropped image of object in test scene. (label index starts from 1.)

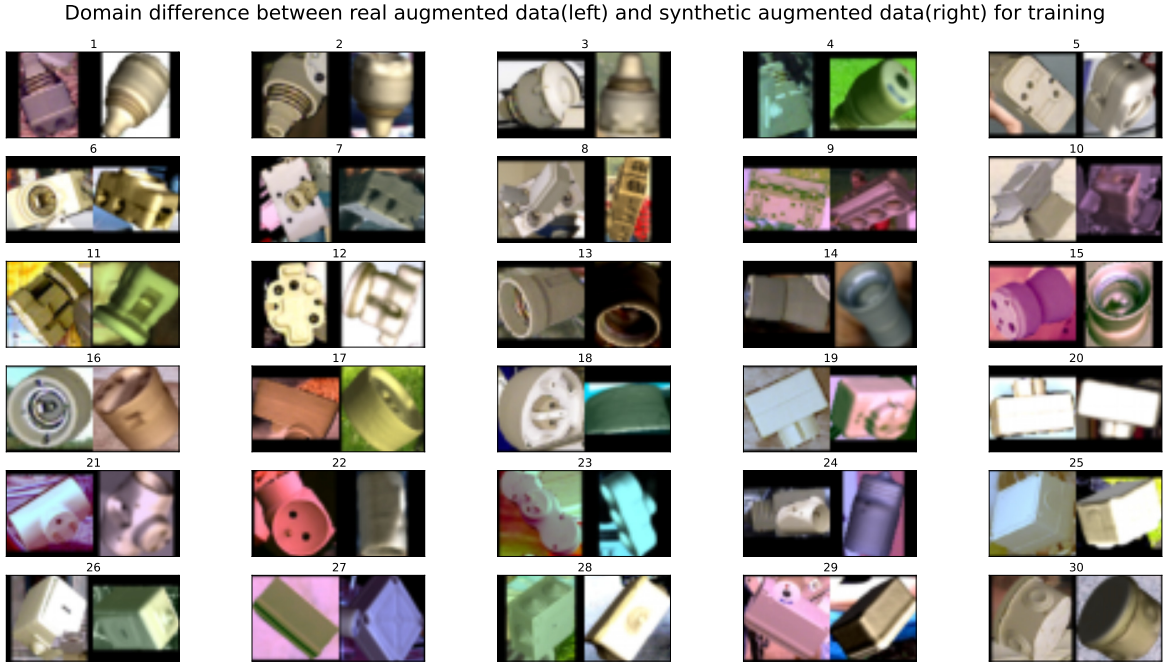


Figure 4.4: Example of object of each category with augmentation. In each thumbnail, the left image is real augmented object and the right image is synthetic augmented object.(label index starts from 1.)

4.1.2 Uncertainty measure

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 measures for uncertainty of prediction, which are introduced in the following. We define \mathbf{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.

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

where $\mathbf{conf} \in [0, 1]$, $c \in \mathcal{L} = \{0, \dots, |\mathcal{L}| - 1\}$, and \mathcal{L} 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]:

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

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

Mutual information between the prediction \mathbf{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:

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

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 metric

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 wr.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 $\mathcal{S}(p(y|x), (y|x))$ that evaluates the quality of predictive distribution $p(y|x)$ relative to an event $y|x \sim q(y|x)$ where $q(y|x)$ represents the true distribution over $(y|x)$. Consequently, the expected scoring rule is:

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

$\mathcal{S}(p, q)$ is proper if $\mathcal{S}(p, q) \leq \mathcal{S}(q, q)$, with equality holds if and only if $p(y|x) = q(y|x)$. 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)** 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.

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

where c_i is the ground truth label for x_i .

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

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

where y_i^{gt} is the one-hot encoding ground truth label for 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

In this part, we evaluate performance of accuracy and uncertainty estimation as well as their summarized performance of different inference techniques for BNN on WRGB and UniHB dataset. Before that, the shorthands of different techniques should be explained. We use "ori" for original ResNet50, "cdp" for ResNet50 with concrete dropout, "mdp" for ResNet50 with multiple dropout without specifications.

In the first experiment, we evaluate them on a relatively simple task, instance recognition. However, our main focus is in the second experiment, where different techniques including dropout, Laplace approximation and ensembles are evaluated and compared on a more difficult task, category recognition. Besides comparing them, we also conduct an ablation study to investigate how feature extractor influences the predictive uncertainty.

4.2.1 Experiments I: uncertainty estimation on instance recognition

In this experiment, we want to evaluate our model on a relatively easy task which is instance recognition at first. 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 $\sim 71.0K$. 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 $\sim 35.6K$ and that of OOD dataset is around $\sim 17.8K$.

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 Fig.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 plot the calibration curve of each model evaluated on both test set and OOD dataset in Fig.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 original ResNet50 is highly overconfident. And the concrete dropout and multi-dropout can mitigate the undesired problem significantly.

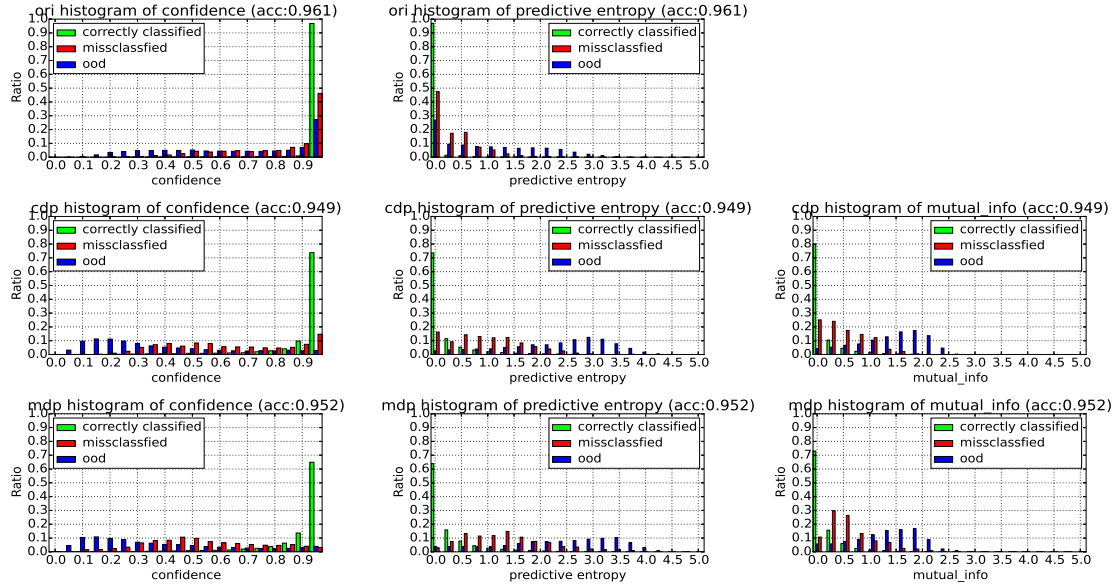


Figure 4.5: Uncertainty(confidence, predictive entropy, mutual information) histograms of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout.

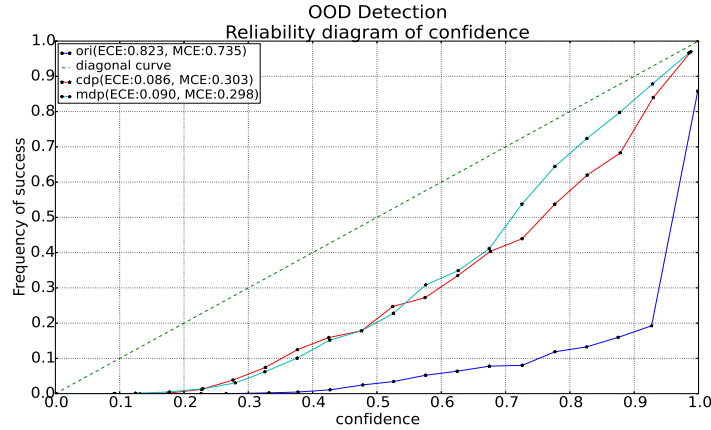


Figure 4.6: Calibration curves of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout evaluated on both test set and OOD dataset.

4.2.2 Experiments II: uncertainty estimation on category recognition

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. Fig.4.1) 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 $\sim 160.9K$, to train our model (in which we split off training set for validation set with ratio 2:8). 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 $\sim 17.1K$. We test performance of uncertainty estimation on this adaptation dataset. The images of 45° , whose size is around $\sim 8.6K$, 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. Fig.4.2) whose size is around $\sim 6.0K$.

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.

In the Fig.4.8, on the left there are 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.

In the Fig.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.

When it comes to the separability metrics in Fig.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.

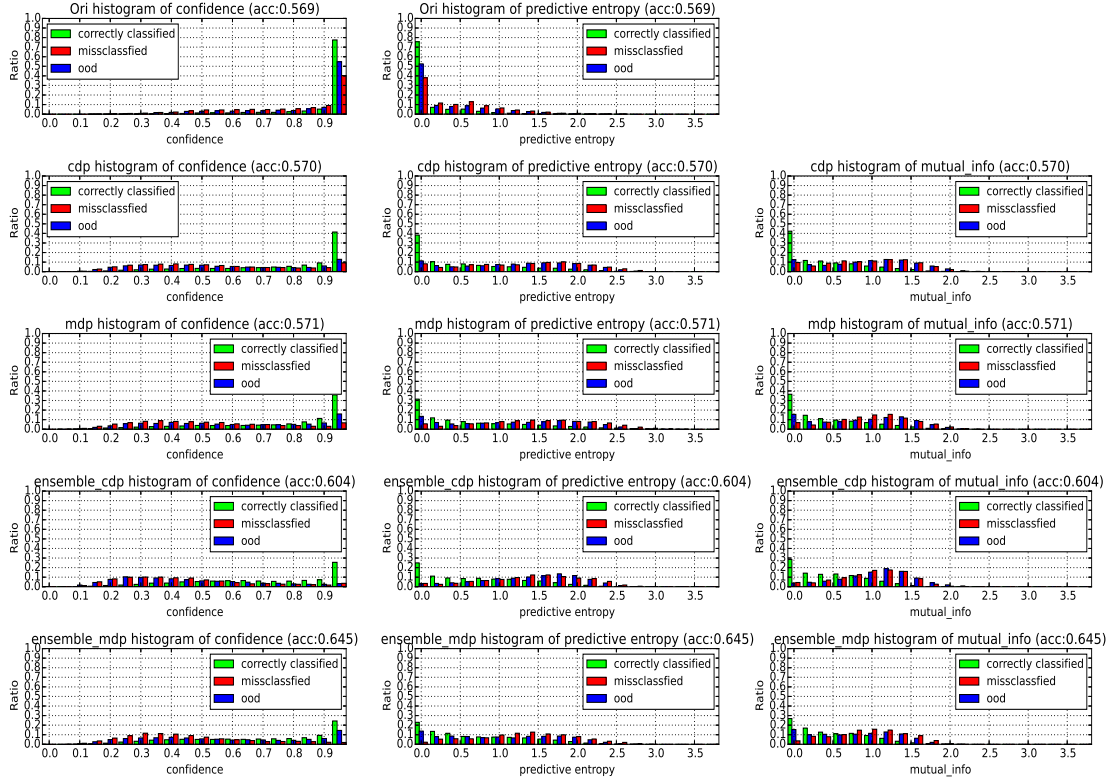


Figure 4.7: Uncertainty histograms of confidence, predictive entropy, mutual information for ori, cdp, mdp, ensemble of cdp, ensemble of mdp in one of three runs.

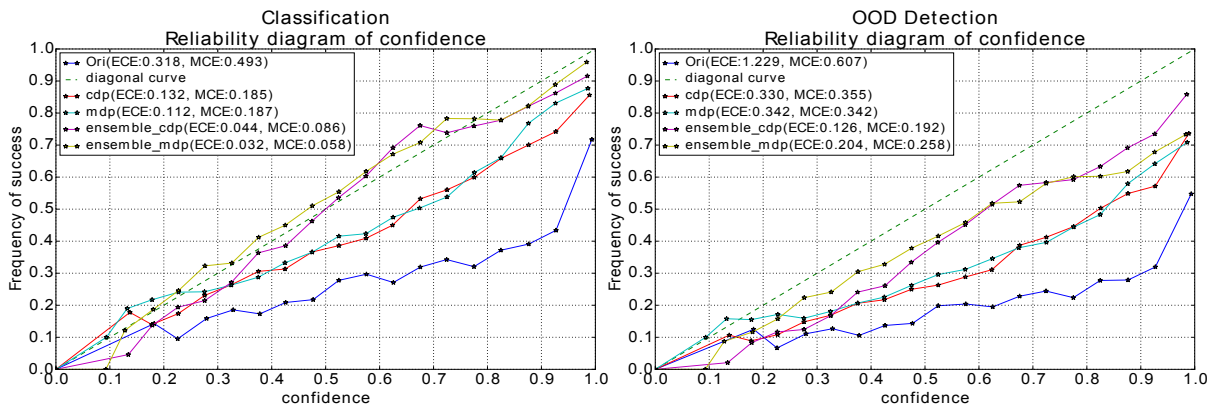


Figure 4.8: Calibration curve for ori, cdp, mdp, ensemble of cdp, ensemble of mdp in one of three runs.

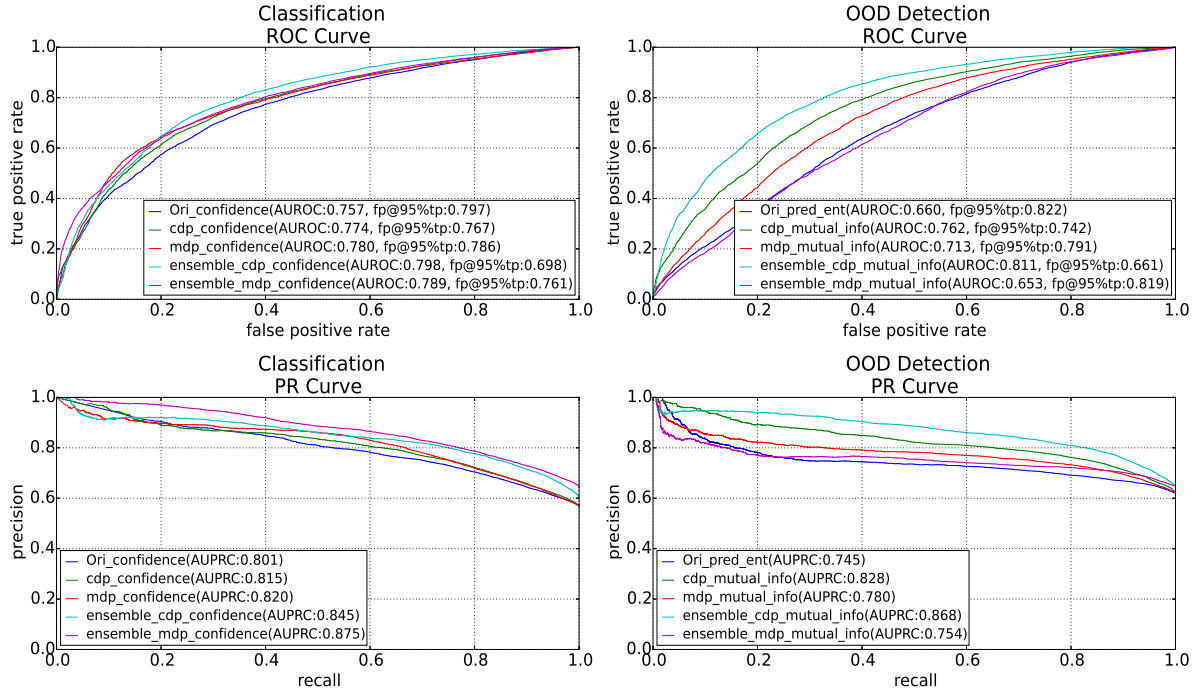


Figure 4.9: ROC and PR curve of ori, cdp, mdp, ensemble of cdp, ensemble in one of three runs.

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 ± 0.008	0.722 ± 0.019	3.242 ± 0.340
cdp	0.577 ± 0.008	0.594 ± 0.013	2.088 ± 0.181
mdp	0.599 ± 0.023	0.566 ± 0.020	1.940 ± 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 ± 0.016 / 0.633 ± 0.065	0.461 ± 0.027 / 0.362 ± 0.025	0.750 ± 0.007 / 0.664 ± 0.011	0.802 ± 0.008 / 0.751 ± 0.018
cdp	0.124 ± 0.023 / 0.288 ± 0.048	0.206 ± 0.015 / 0.374 ± 0.018	0.775 ± 0.008 / 0.783 ± 0.022	0.825 ± 0.007 / 0.850 ± 0.022
mdp	0.114 ± 0.012 / 0.383 ± 0.046	0.199 ± 0.016 / 0.367 ± 0.023	0.780 ± 0.011 / 0.709 ± 0.004	0.838 ± 0.013 / 0.788 ± 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 2.37 based on grid search on the validation set.

Different visual metrics such histograms, calibration curves, ROC curve and PR curve of cdp, mdp and their Laplace approximation versions are showed in the following. 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.

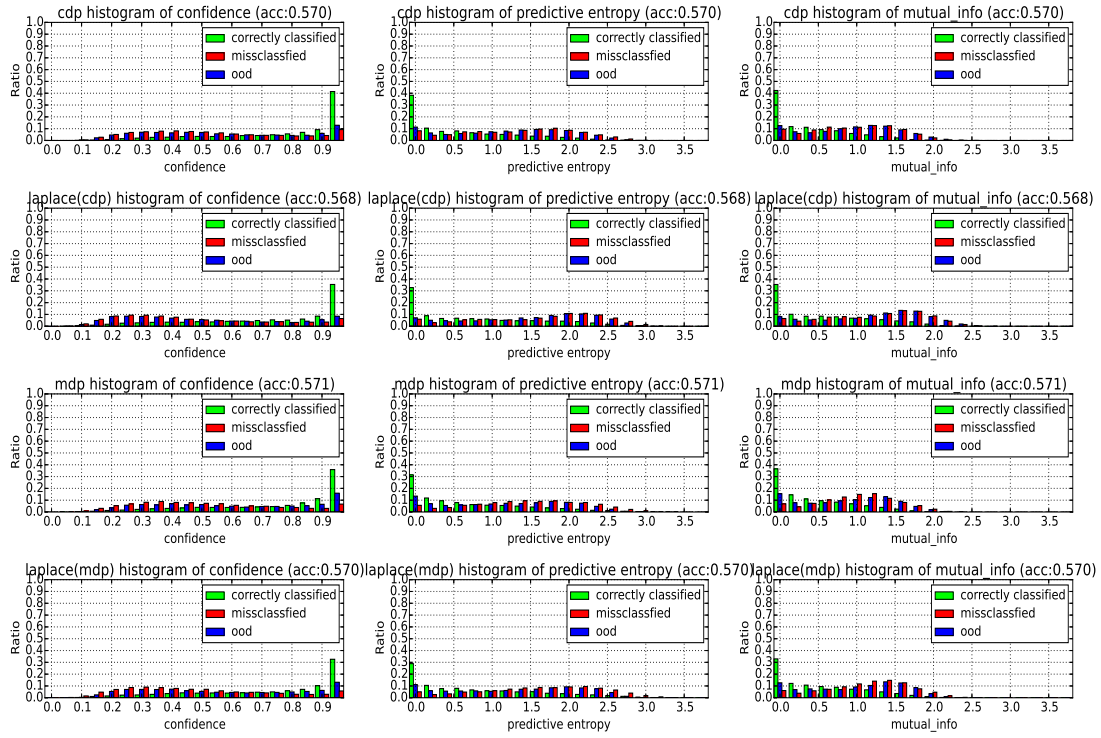


Figure 4.10: Uncertainty histogram of cdp, mdp and their Laplace approximation versions with confidence, predictive entropy and mutual information as uncertainty measure in one of three runs.

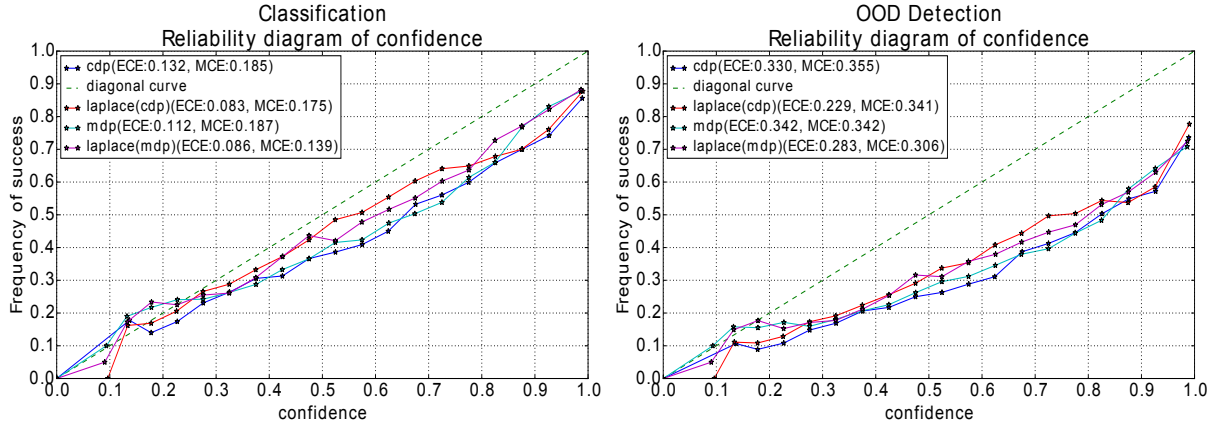


Figure 4.11: Calibration curve of cdp, mdp and their Laplace approximation versions in one of three runs.

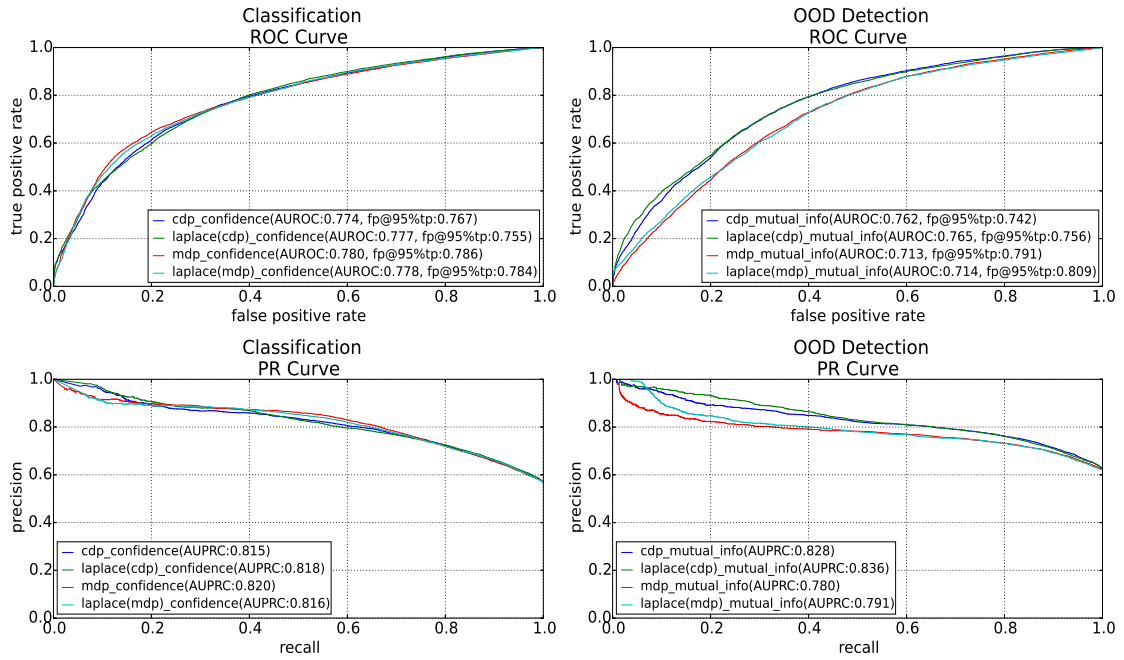


Figure 4.12: ROC and PR curve of cdp, mdp and their Laplace approximation versions in one of three runs.

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 an ablation study to check the influence of aleatoric uncertainty which comes from mainly the feature extractor part. To this end, we firstly make assumption that the feature extractor part has a dominant effect on the aleatoric uncertainty, because it's deterministic compared with the following MLP classifier.

Considering this, we freeze the parameters of feature extractor trying to keep the dominant effect on aleatoric uncertainty fixed for different approaches 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 here. 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.

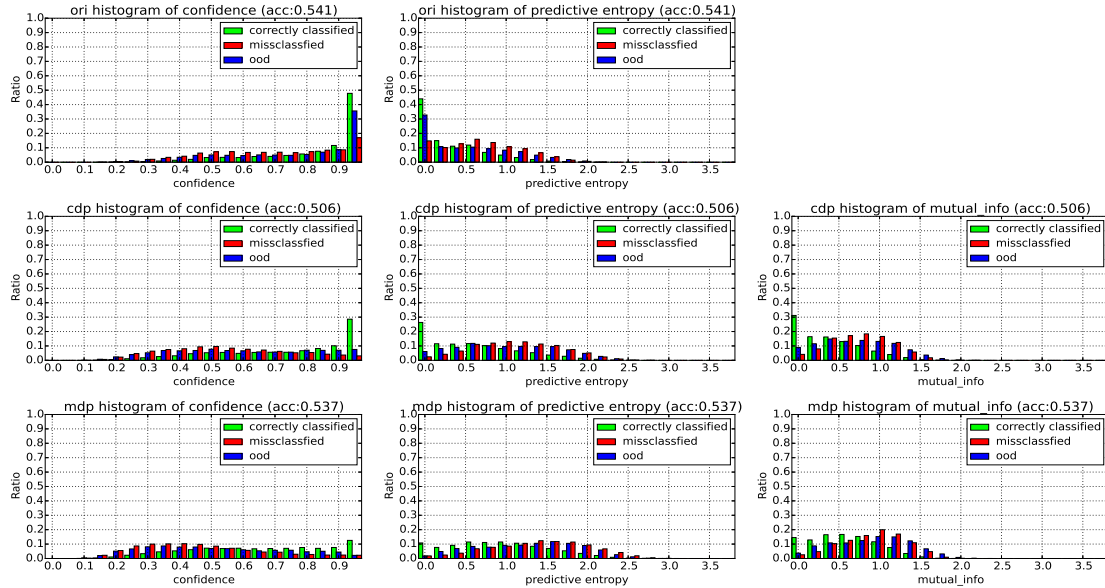


Figure 4.13: Uncertainty histograms of ori, cdp, mdp trained with frozen features in one of three runs.

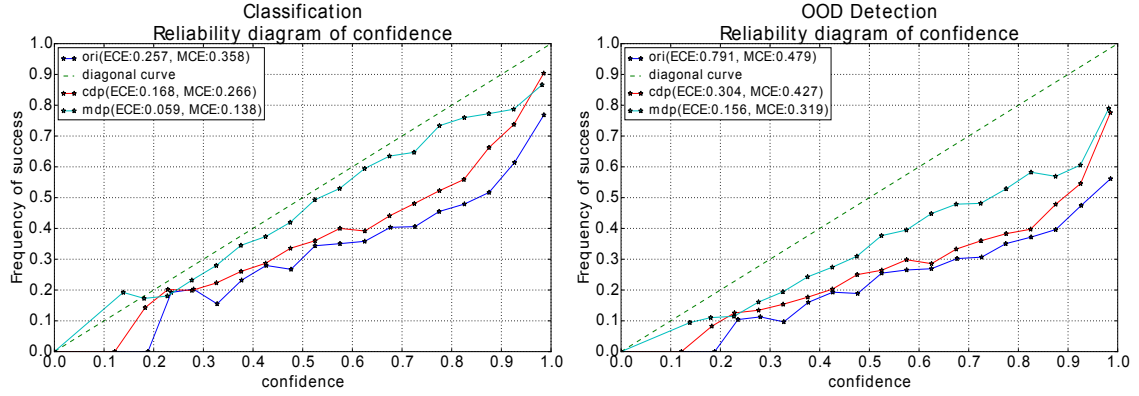


Figure 4.14: Calibration curves of ori, cdp, mdp trained with frozen features one of three runs.

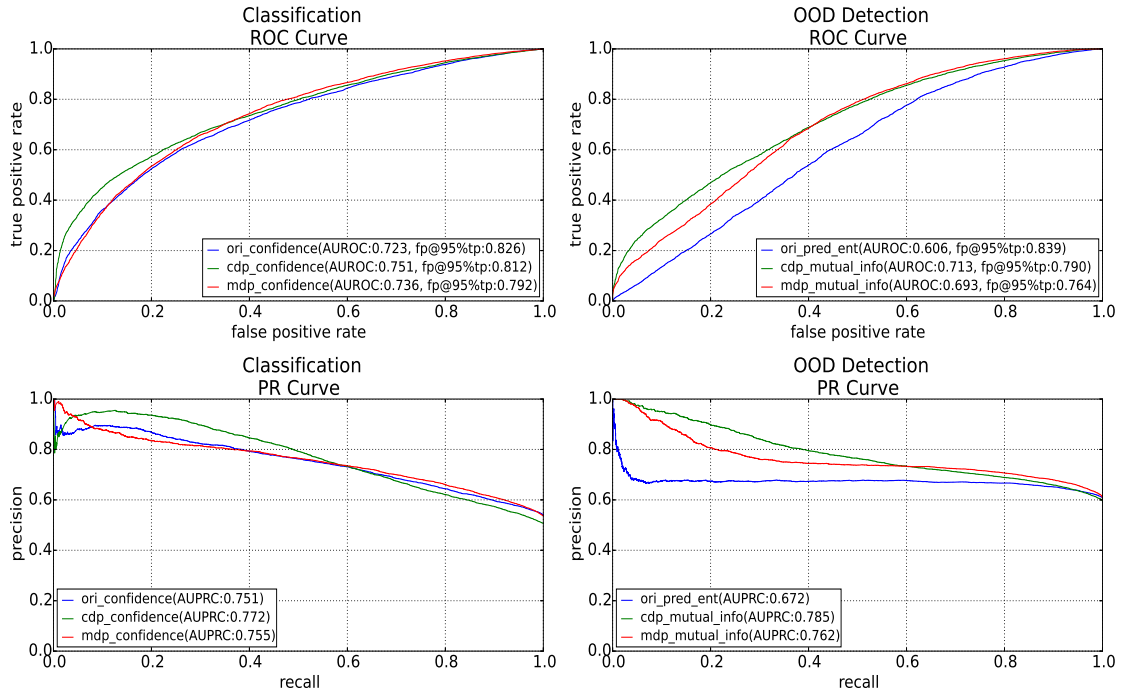


Figure 4.15: ROC and PR curves of ori, cdp, mdp trained with frozen features one of three runs.

Table 4.3: Quantitative results averaged from 3 random seeds.

	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

In experiments of automatic labeling, 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 (cf. 3.6). The automatically labeled data is collected based on improved uncertainty estimation. Because ResNet50 with concrete dropout shows more stable results than other version, we use it to improve uncertainty estimation in the following experiments and call it **BNN** without specifications.

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

Then we try to extend this way to another dataset, T-LESS dataset, to test its generalization ability in second experiment. In the second experiment, the situation is similar to the first experiment, where a classifier is firstly trained on synthetic dataset and then deployed in a real world environment which is simulated by the real objects in training set. We employ data augmentation to address the problem of imbalance in fine-tuning classifier found in the first experiment.

4.3.1 Experiment I: evaluation on WRGBD and UniHB dataset

The first experiment is performed on the WRGBD and UniHB dataset. As is mentioned in experiment II of uncertainty estimation experiments, we train BNN with entire WRGBD dataset and evaluate performance of uncertainty estimation on objects of elevation 30° and 60° in UniHB dataset, which we call **adaptation dataset**. The size of adaptation dataset in this experiment is around $\sim 17.1K$. The objects of 45° with size $\sim 8.6K$ 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.

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 misclassifications based on the results of previous experiments.

Automatic Labeling: 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 10 data samples in each category. Then we can investigate how two factors may influence the performance of fine-tuning. These two factors are:

- the balance in number of data sample of each class
- the amount of information of data samples

In the following table 4.3.1, we describe different settings and use capital letter to denote them.

Here we firstly illustrate the steps of **selecting automatically labeled data** in this experiment which is a little 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 this context means that we randomly choose the data of each class with number computed based on the difference between current number of data samples of each class and the average number of data samples of each class. We assume that we have manually labeled all the predictions after automatic labeling and conduct a *proof-of-concept* experiment with **as balanced as possible** dataset in F and G.

Table 4.4: Results of fine-tuned network with 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%

Fine-tuning: After collecting the dataset with different settings, we show the results of BNN fine-tuned with those dataset on the final test set which contains only objects of 45° from UniHB data in the table 4.3.1. We can draw the following observations from this table:

- 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 least 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 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.

4.3.2 Experiment II: evaluation on T-LESS dataset

In this experiment, original ResNet50 and BNN are trained on synthetic dataset introduced in section 4.1.1. Then automatically labeled data is collected by testing the trained model on real single objects (original training set of T-LESS) which we call **adaptation dataset**. The difference between the strategy of selecting data in experiment I is that threshold used for selecting data based on uncertainty is chosen on validation set. The accuracy of predictions with uncertainty smaller than this threshold is set as 95%. Validation set is split off from training set with ratio 2:8. The reason why we use this way to choose threshold is that it's more practical and closer to the real world applications. In this experiment predictive entropy is selected as uncertainty measure, because it shows better performance in term of separability metrics on validation set.

The improved uncertainty estimation from BNN plays an important role here, because it not only provides a reliable uncertainty estimation, but also increases the separability between correct predictions and false predictions, which is more useful in this task.

Automatic labeling: We show the comparisons of uncertainty estimation in term of different evaluation metrics with BNN and original ResNet50 in the following figures. As can be seen in the uncertainty histogram (cf.4.16), the mode of distribution of uncertainty estimation between correct prediction and miss-classification is further from each other when testing with BNN than with original ResNet50. In Fig.4.17, BNN expresses better performance in terms of both calibration and separability metrics.

With the same procedure of selecting automatically labeled data, the size of dataset is around $\sim 1.05K$ with only 93% accuracy using original ResNet50, but around $\sim 1.6K$ with

96% accuracy using BNN. Because original ResNet50 produces less automatically labeled data with lower quality, the dataset for fine-tuning is collected with BNN.

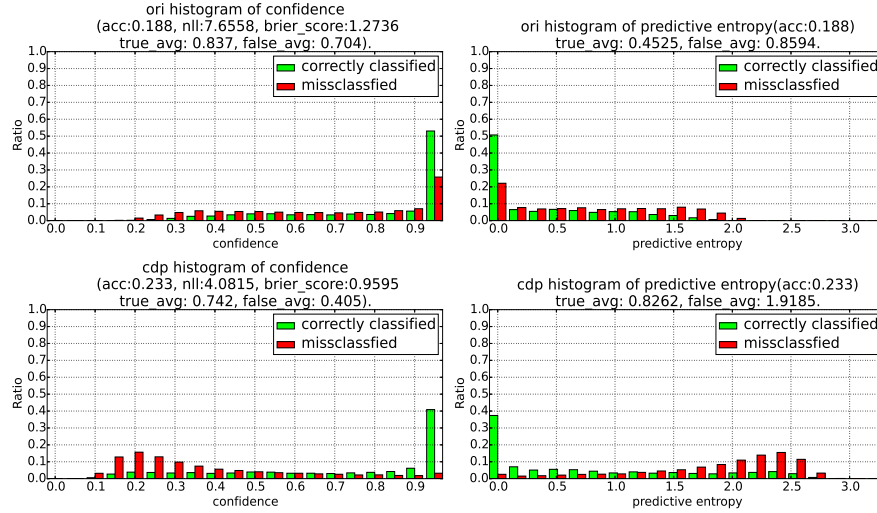


Figure 4.16: Uncertainty histograms of original ResNet50 (top) and BNN (down).

However, as found in experiment I, the problem of imbalance in number of data samples of each class still appears on this dataset and has a large influence on the performance. The horizontal histogram in Fig.4.18 reveals this problem, in which most of confident predictions cluster in few classes which are class 6, 7, 28 and so on. To mitigate this problem, two simple ways are adopted here:

- to manually label data uniformly with size 3% of the adaptation set (around $\sim 0.9K$ data samples).
- to employ augmentations described in section 4.1.1 to the objects to balance the dataset as much as possible.

The balance of dataset with/without manually labeled data after augmentations are showed in Fig.4.18. By using augmentations, datasets are more balanced than before. But the balance before augmentations has significant effect on the performance, which is discussed in the following.

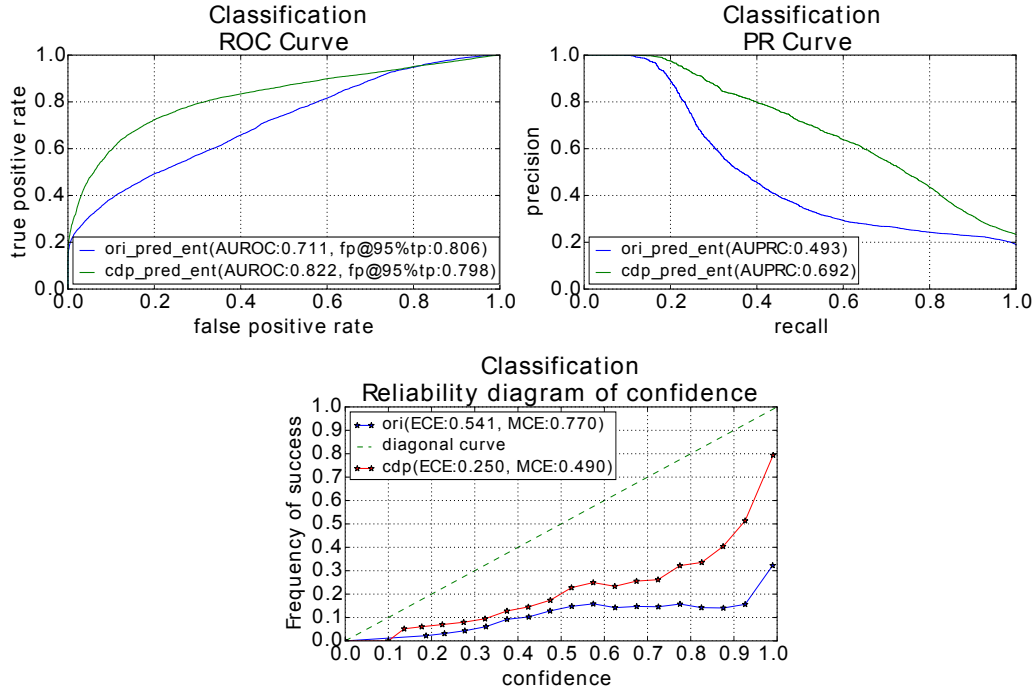


Figure 4.17: Calibration curves, ROC curves, PR curves of original ResNet50 and BNN.

Fine-tuning: BNNs trained with only synthetic data previously are fine-tuned on dataset collected with different strategies. Table 4.3.2 gives the results of BNN trained with different dataset. They are all tested on the test set of T-LESS dataset. By comparing A and B, the performance drop caused by the domain gap between synthetic objects and real ones is quite large, around 35%. By fine-tuning model on dataset collected based on uncertainty estimation, the accuracy increases a lot. With more balanced dataset, the accuracy even exceed the model trained with augmented entire real dataset by around 3%. By comparing C and D, it's obvious that the balance of data samples in each class before augmentation is important to obtain more accuracy gain. The reason behind that is the diversity of data samples before augmentation. Although augmentation can help the model generalize better by increasing the number of data samples, if the distribution of training set is poorly represented by few training samples, the way does not help much. This can also be observed in comparison of the corresponding confusion matrices in Fig.4.19, even though the number of data samples in some classes are similar after augmentation, their accuracies differ a lot (35% vs. 65% in class 0).

Table 4.5: Results of BNN trained with different dataset (size of dataset before augmentations is showed in the bracket).

Setting of dataset	Accuracy
A: augmented, synthetic dataset	34.91%
B: augmented, entire real dataset ($\sim 30K$)	69.58%
C: fine-tune A with augmented, only automatically labeled real dataset ($\sim 1.6K$)	44.99%
D: fine-tune A with augmented, automatically labeled and 3% manually labeled real dataset ($\sim 2.5K$)	72.48%

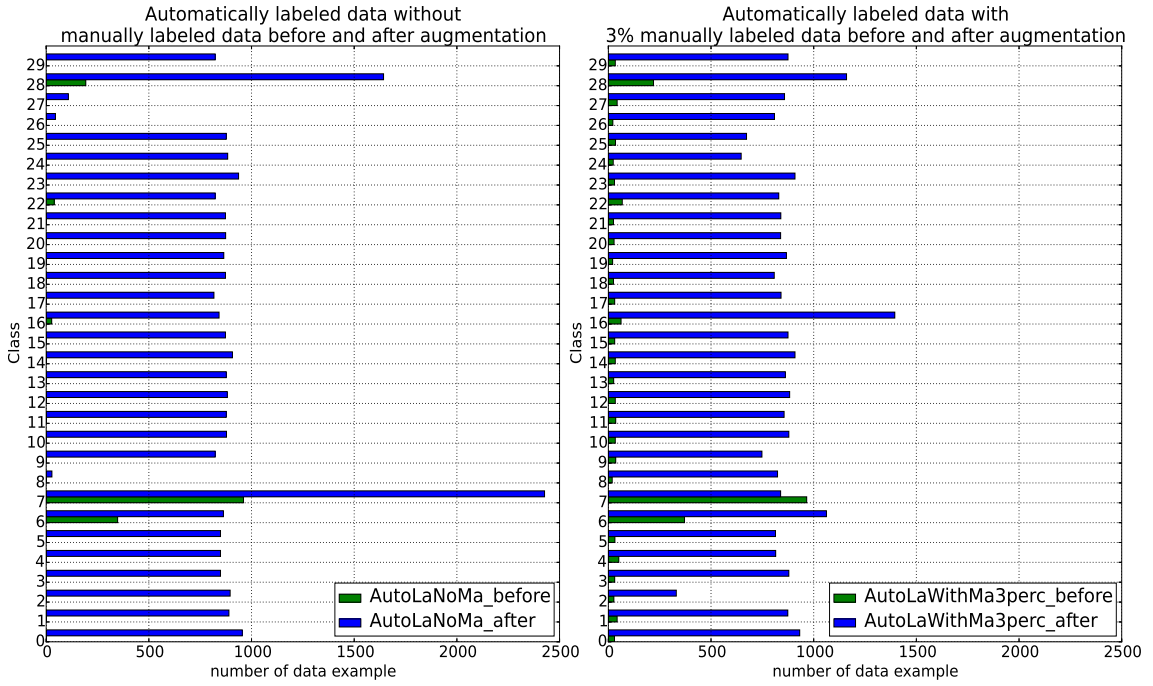


Figure 4.18: Horizontal histograms of automatically labeled data with/without 3% manually labeled data before/after augmentations.

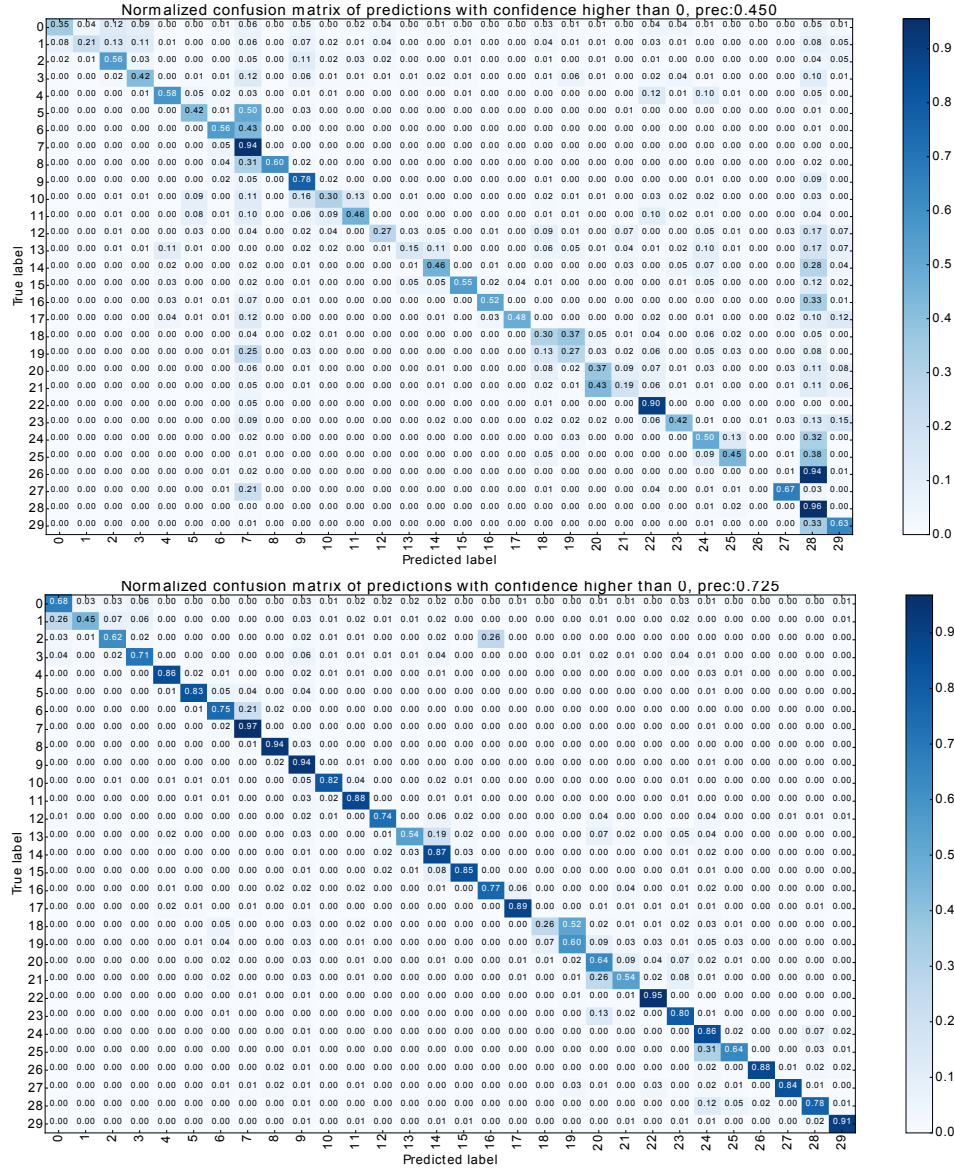


Figure 4.19: Confusion matrices of BNN fine-tuned with dataset with/without small portion of manually labeled data.

4.4 Context-based improvement experiments

In this section, experiments are performed to investigate the effect of combining improved uncertainty estimation with a probabilistic graphical model in down-stream task. In detail, we employ conditional random field (CRF) to capture the dependencies between objects by exploiting contextual information in the scene(cf. Fig.3.3). The intuition behind this idea is that, if the predictive distribution from BNN is more reliable and can represent the inherent randomness better. Then this bit of information can be utilized by CRF

for further performance gain. Similar to previous experiments, ResNet50 with concrete dropout is used to improve uncertainty estimation in the following experiments and we call it **BNN** without specifications.

Firstly a simpler experiment is performed on subset of T-LESS dataset to show that CRF can capture the dependencies and improve the accuracy, which can be boosted with better uncertainty estimation from BNN. Secondly, the same idea is tested on the entire T-LESS dataset to show the generalization ability of this approach.

4.4.1 Experiment I: evaluation on subset of T-LESS dataset

In this experiment, original ResNet50 and BNN are trained on the real single objects with black background (original training set of T-LESS dataset). Then a subset of test scenes of original T-LESS dataset are used for training and testing CRF. In detail, scene 1, 2, 3, 4 ($\sim 8K$ objects and $\sim 2K$ scenes) are used for training CRF, and scene 5, 6, 7, 8 ($\sim 12K$ objects and $\sim 2K$ scenes) for testing. To note that **no augmentations** are applied in this experiment. To note that, based on the prior knowledge about the testing set, that no objects repeat in the test scenes, we set the diagonal entries of binary co-occurrence matrix with -1 to penalize the outcome with two same objects.

In the table 4.6, results of CRF trained with two different unary features from original ResNet50 and BNN are given. As can be seen, CRF trained and tested with unary feature from BNN can yield more improvement(13.17%) compared with that from deterministic network (12.42%).

Table 4.6: Results of CRF trained with unary feature from different versions of ResNet50

	accuracy with unary potential	accuracy with unary and pairwise potentials
unary feature from ori	48.79%	61.21%
unary feature from BNN	51.94%	65.11%

Additionally, because the weights in CRF model act as coefficients for two kinds of features, we can check the magnitude of these two weights for their contributions to overall performance. In the table 4.7, the weight of unary feature θ_u from BNN is larger than the one of edge feature θ_p . The situation is inverse in unary feature from original ResNet50. It's evident that the unary feature from BNN contributes more than the pairwise one while the one from original ResNet50 contributes less. It also shows that the predictive distribution from BNN contains more information than its counterpart from original ResNet50. More importantly, this information can be utilized to improve the performance which is shown in the table 4.6.

Table 4.7: Learned weights of CRF trained with unary feature from different versions of ResNet50

	node weight θ_u	edge weight θ_p
unary feature from ori	3.744	4.213
unary feature from BNN	4.635	3.882

4.4.2 Experiment II: evaluation on entire T-LESS dataset

In this experiment, approach of combination of BNN and CRF is tested on the whole T-LESS dataset. Based on the results obtained in subsection 4.3.2, the accuracy should be improved further by exploiting contextual information with CRF (cf. Fig.3.5). In details, CRF are trained with unary features extracted from BNN before fine-tuning with size around ~ 24.5 objects and $\sim 5K$ scenes (cf. Fig.3.6). In testing, the whole testing set of original T-LESS dataset with size $\sim 70K$ objects and $\sim 10K$ scenes, whose unary features are extracted from fine-tuned BNN.

Because the combinations of the whole test set are more complex than that of scene 5, 6, 7, 8 in experiment I, we try different values on the diagonal with which different prior knowledge is encoded. There are many objects occurring repeatedly in other scenes of test set besides scene 5, 6, 7, 8. Therefore there is a trade-off between different kind of prior knowledge. Filling diagonal entries with 1 indicates that co-occurrence of same objects is more likely to happen and it suits better on the whole test set. As can be seen in the table 4.4.2, CRF trained with pairwise matrix with 1 on diagonal has most improvement (4%) when comparing with matrices with other diagonal values.

Table 4.8: Results of CRF with different type of potentials

Type of potential	Accuracy
Unary	72.48%
Unary and pairwise(-1s on diagonal)	73.40%
Unary and pairwise(0s on diagonal)	74.64%
Unary and pairwise(1s on diagonal)	76.50%

Chapter 5

Summary and Conclusion

5.1 Summary

In this work, different approaches for improvement of a classification system related to uncertainty have been investigated.

Firstly, in order to enable classifier to adapt itself in test environment which has domain gap with its training set, dataset for fine-tuning should be collected with as little human efforts as possible. We propose to gather automatically labeled data based on uncertainty estimation of the classifier. Considering the overconfident behavior of deep learning classifier, we employ Bayesian neural network to improve the uncertainty estimation. For easy compatibility with existing neural network architectures, dropout variational inference and scalable Laplace approximation chosen for this task. In addition to the naive version, advanced variants of them such as concrete dropout, multiple dropout and ensemble of them are taken into account. We modify ResNet50 to incorporate these techniques and keep the ability of extracting powerful features at the same time. Extensive experiments on WRGBD and UniHb dataset are performed to evaluate the performance of uncertainty estimation of different approaches in terms of different metrics including accuracy, proper scoring rules, calibration metrics, separability metrics. The results show that Bayesian neural network can improve the performance of both accuracy and calibration as well as uncertainty estimation when compared with the original ResNet50.

Secondly, the improved uncertainty estimation is used for collecting automatically labeled dataset, which is used for fine-tuning the model trained on a easily obtainable dataset. The goal is to obtain a more accurate domain-specific classifier with as little human efforts as possible. Experiments of this part are performed on both WRGBD/UniHB dataset and original T-LESS/synthetic T-LESS dataset. During experiments, we find that the balance of data samples in each class has a significant influence on the performance after fine-tuning. Augmentations are employed to mitigate this problem and the experiments results show that, a more accurate domain-specific can be obtained with much less human efforts,

while its performance can reach or even exceed the performance of classifier obtained with much more human efforts.

Thirdly, the improved uncertainty estimation should be utilized further with proper probabilistic model. On the other hand, uncertainty induced by appearances similarity should be handled by modeling the contextual information between objects. On account of these two points, conditional random field is employed to further improve the performance by utilizing more information from out put of Bayesian neural network and capturing the contextual relationship between objects in the same scene. For the end, experiments are performed to validate this idea but only on T-LESS dataset because designing pairwise feature in WRGBD dataset is non-trivial. The experiments results show that CRF can improve the accuracy further by not only making use of the information from better uncertainty estimation but exploiting the contextual relationship between objects.

Although experiments results show that these approaches can achieve better performance compared with original version, the problems occurring when constructing approaches and performing experiments should require more attention and efforts towards solving them. There are several main problems listed in the following:

- inefficiency in testing with Bayesian neural network: testing each input requires sampling the posterior many times. Different approaches are already available such as network[HVD15] distillation and parallel computation.
- Uncertainty estimation should be defined more clearly and explored with more efforts. For example, difference of separability metrics and calibration metrics are used to evaluate uncertainty estimation in literatures[HG16][GPSW17]. However, they should account for different ability which should not be summaries with one word. As can be seen in the experiments, well calibrated model may have similar or even worse separability between correct predictions versus miss-classifications or OOD predictions than the badly-calibrated model.
- Other approaches to improve uncertainty estimation should be considered. Bayesian neural network is a principal way for this purpose. But the computation behind that is too high. It's like using a generative model to solve a problem which can be solved by discriminative model. Maybe there is more efficient way to achieve this goal such as different post-processing methods [GPSW17].
- the way to choose highly confident predictions as automatically labeled data for continuous learning sounds practical and meaningful. However, from perspective of active learning, the confident predictions contain little information required by the classifier. To train a classifier with better performance requires training set with more information. To label predictions with more information (low confidence) sounds more meaningful. Therefore a better technique to combine them should be taken into account.
- modeling dependencies between random variables is a promising and natural way to

improve the results because the i.i.d. assumption does not hold all the time. However, employing probabilistic graphical model requires hand-crafted edge feature which is non-trivial to have and design. Therefore, graph neural network which can learn the edge feature based on data is promising to work on.

5.2 Conclusion

- Bayesian neural network can improve uncertainty estimation in terms of different evaluation metrics including accuracy, proper scoring rule, calibration and separability between correct prediction versus miss-classification or OOD prediction, while ensemble of them can achieve even better performance.
- Multiple dropout can work better without considering OOD data, while concrete dropout can offer more stable performance on separating correct predictions and miss-classifications.
- Scalable Laplace approximation with Kronecker-factored approximation can achieve similar performance as concrete dropout and multiple dropout, indicating that the performance is highly related to the point estimate obtained during training.
- Manual efforts in collecting dataset for fine-tuning to enable classifier to adapt itself in test environment can be reduced with the help of better uncertainty estimation. Nevertheless, imbalance of collected dataset has a negative effect on fine-tuning, which can be mitigated by data augmentation and adding little manual labeling.
- Conditional random field can improve performance further by utilizing information brought by better uncertainty estimation and exploiting contextual information between objects in the same scene.

5.3 Future work

There are interesting work related to the approaches used in this work:

- Uncertainty distribution of different view points of the object can be checked. On the one hand, to understand the behavior of neural network better with uncertainty. On the other hand, if the uncertainty can represent actual discriminability of different view points, this information is useful in down-stream tasks such optimal view point planing for grasping and so on.
- As mentioned in summary, inefficiency in testing is an important problem. It would be interesting to try distillation method to decrease runtime in testing.
- Although in this work, dropout variational inference and Laplace approximation are evaluated. The approximation distribution of them seem restrictive and simple.

Interesting future work would be to try more complex and expressive approximate posterior distribution such as normalizing flow .

- In the experiments of evaluating uncertainty estimation, multiple dropout can achieve better performance without considering OOD data. Ablation study has also showed that multiple dropout can have better performance. As stated in [KSW15], covariances between instances increase variances of estimating derivatives of variational parameters. One way to improve performance of multiple dropout would be to decrease the variance of estimating variational parameters by sampling masks for each instance separately because there are more parameters than concrete dropout.

Appendix A

Appendix

A.1 Fisher information

Here, we consider the case that parameter of likelihood function is scalar. This can easily be extended to the case that θ is a vector. Fisher information $\mathcal{I}(\theta)$ is defined as follows:

$$\mathcal{I}(\theta) = \mathbb{E}\left[\left(\frac{\partial}{\partial\theta} \log f(\mathbf{X}; \theta)\right)^2\right] \quad (\text{A.1})$$

Where $f(\mathbf{X}; \theta)$ is probability density function or mass function of random variable \mathbf{X} . We need to prove that $\mathbb{E}\left[\left(\frac{\partial}{\partial\theta} \log f(\mathbf{X}; \theta)\right)^2\right] = \mathbb{E}\left[-\frac{\partial^2}{\partial\theta^2} \log f(\mathbf{X}; \theta)\right]$. When we suppose that θ is the true parameter. Based on the fact that $f(\mathbf{X}; \theta)$ is density, we can have:

$$\int f(x; \theta) dx = 1$$

Differentiating the above with respect to θ gives:

$$\frac{\partial}{\partial\theta} \int f(x; \theta) dx = 0$$

After moving differentiating operation into integration and making use of derivative of logarithm, we get:

$$\int \frac{\partial \log f(x; \theta)}{\partial\theta} f(x; \theta) dx = 0$$

Differentiating again w.r.t. θ and taking the derivative inside gives:

$$\begin{aligned}
& \int \frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} f(x; \theta) dx + \int \frac{\partial \log f(x; \theta)}{\partial \theta} \frac{\partial f(x; \theta)}{\partial \theta} dx = 0 \\
\Rightarrow & \int \frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} f(x; \theta) dx + \int \frac{\partial \log f(x; \theta)}{\partial \theta} \frac{\partial f(x; \theta)}{\partial \theta} \frac{1}{f(x; \theta)} f(x; \theta) dx = 0 \\
\Rightarrow & \int \frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} f(x; \theta) dx + \int \frac{\partial \log f(x; \theta)}{\partial \theta} \frac{\partial \log f(x; \theta)}{\partial \theta} f(x; \theta) dx = 0 \\
\Rightarrow & -\mathbb{E}\left[\frac{\partial^2 \log f(x; \theta)}{\partial \theta^2}\right] = \mathbb{E}\left[\frac{\partial \log f(x; \theta)}{\partial \theta} \frac{\partial \log f(x; \theta)}{\partial \theta}\right]
\end{aligned}$$

which gives us the required result. □

A.2 Derivation of Re-parameterization Trick

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

$$\begin{aligned}
I'(\theta) &= \frac{\partial}{\partial \theta} \int f(x) p_\theta(x) dx \\
&= \frac{\partial}{\partial \theta} \int f(x) \left(\int p_\theta(x|\epsilon) p(\epsilon) d\epsilon \right) dx \\
&= \frac{\partial}{\partial \theta} \int \left(\int f(x) p_\theta(x|\epsilon) dx \right) p(\epsilon) d\epsilon \\
&= \frac{\partial}{\partial \theta} \int \left(\int f(x) \delta(x - g(\theta, \epsilon)) dx \right) p(\epsilon) d\epsilon \\
&= \frac{\partial}{\partial \theta} \int f(g(\epsilon, \theta)) p(\epsilon) d\epsilon \\
&= \int \frac{\partial}{\partial \theta} f(g(\epsilon, \theta)) p(\epsilon) d\epsilon \\
&= \int f'(g(\epsilon, \theta)) \frac{\partial}{\partial \theta} g(\theta, \epsilon) p(\epsilon) d\epsilon \\
&= \mathbb{E}_{p(\epsilon)} \left[f'(g(\epsilon, \theta)) \frac{\partial}{\partial \theta} g(\theta, \epsilon) \right]
\end{aligned} \tag{A.2}$$

A.3 Derivation of Hessian of Log Posterior

$$\begin{aligned}
\log p(\omega|Y, X) &\approx \log p(\omega^*|Y, X) + \\
&\quad (\omega - \omega^*)^T \frac{\partial \log p(\omega|Y, X)}{\partial \omega} + \\
&\quad \frac{1}{2} (\omega - \omega^*)^T \frac{\partial^2 \log p(\omega|Y, X)}{\partial \omega^2} (\omega - \omega^*) \\
&= \log p(\omega^*|Y, X) - \frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*)
\end{aligned} \tag{A.3}$$

where

List of Figures

2.1	Differences between parameter estimation of the deterministic neural network and the BNN. The solid point denotes deterministic variable, and the circle denotes random variable, while the shaded circle denotes observed random variable. The plate notation denotes \mathbf{N} observed data pairs (\mathbf{x}, \mathbf{y}) .	9
2.2	Schematic illustration of dropout in a three-layer fully connected neural network. (a) the standard neural network with all neurons switched on. (b) after applying dropout only a subset of neurons are used[SHK ⁺ 14].	11
2.3	An example of two layer neural network with dropout, a Bernoulli random variable is imposed on each unit of input layer and hidden layer.	12
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($\mathbf{0}, \mathbf{1}$) [MMT16].	20
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($\mathbf{0}, \mathbf{1}$)[MMT16].	21
2.6	One sample and average value of 100 samples drawn from continuous approximation of Bernoulli distribution with parameter $\mathbf{p} = [\mathbf{0.1}, \mathbf{0.2}, \dots, \mathbf{1.0}]$ and temperature $\lambda = \mathbf{0.1}$	22
2.7	Simple example of conditional random field.	29
3.1	Changes of keep probability of the first two concrete dropout layers during training.	34
3.2	Changes of keep probability of the last two concrete dropout layers during training.	35
3.3	Different dropout rates for different hidden units in multi-drop.	35
3.4	Modified network architecture of ResNet50.	37
3.5	Combination of BNN and CRF.	38
3.6	Approach for continuous learning.	39

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.	43
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. .	43
4.3	Example of object of each category in original training set and test set. In each thumbnail, the left image is real single object with black background and the right image is cropped image of object in test scene. (label index starts from 1.)	44
4.4	Example of object of each category with augmentation. In each thumbnail, the left image is real augmented object and the right image is synthetic augmented object.(label index starts from 1.)	44
4.5	Uncertainty(confidence, predictive entropy, mutual information) histograms of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout.	50
4.6	Calibration curves of original ResNet50, ResNet50 with concrete dropout and ResNet50 with multi-dropout evaluated on both test set and OOD dataset.	50
4.7	Uncertainty histograms of confidence, predictive entropy, mutual information for ori, cdp, mdp, ensemble of cdp, ensemble of mdp in one of three runs.	53
4.8	Calibration curve for ori, cdp, mdp, ensemble of cdp, ensemble of mdp in one of three runs.	53
4.9	ROC and PR curve of ori, cdp, mdp, ensemble of cdp, ensemble in one of three runs.	54
4.10	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
4.11	Calibration curve of cdp, mdp and their Laplace approximation versions in one of three runs.	57
4.12	ROC and PR curve of cdp, mdp and their Laplace approximation versions in one of three runs.	57
4.13	Uncertainty histograms of ori, cdp, mdp trained with frozen features in one of three runs.	59
4.14	Calibration curves of ori, cdp, mdp trained with frozen features one of three runs.	60
4.15	ROC and PR curves of ori, cdp, mdp trained with frozen features one of three runs.	60
4.16	Uncertainty histograms of original ResNet50 (top) and BNN (down). . . .	65
4.17	Calibration curves, ROC curves, PR curves of original ResNet50 and BNN.	66
4.18	Horizontal histograms of automatically labeled data with/without 3% manually labeled data before/after augmentations.	67

4.19 Confusion matrices of BNN fine-tuned with dataset with/without small portion of manually labeled data.	68
---	----

List of Tables

4.1	Quantitative results of acc, bs, nll, ece, mce, auROC, aupr averaged from 3 different random seeds	55
4.2	Quantitative results of acc, bs, nll, ece, mce, auROC, aupr averaged from 3 different random seeds	58
4.3	Quantitative results averaged from 3 random seeds.	61
4.4	Results of fine-tuned network with different settings	63
4.5	Results of BNN trained with different dataset (size of dataset before augmentations is showed in the bracket).	67
4.6	Results of CRF trained with unary feature from different versions of ResNet50	69
4.7	Learned weights of CRF trained with unary feature from different versions of ResNet50	70
4.8	Results of CRF with different type of potentials	70

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.
- [EVGW⁺] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.

- [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] Tilmann 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] Tilmann 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.
- [HVD15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [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.
- [LLS17] 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. Upgmpp: 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 Yarín 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.