

Computer Vision with Real-World Data

Student	Kirilin Anton Dmitrievich
Student's time capacity	8 ± 3 hours
Simulation time capacity (GPU: 20 GB VRAM)	10 ± 5 hours

General Instruction

Please use Russian or English language when preparing this document, but not their mixture. The exception is for the code in the Appendix section which may be a mixture of operators in English and comments in Russian.

Students are welcome to use any tool that is suitable for preparing high-quality work. However, we ask students to keep in mind two important criteria. First, we expect students to fully describe their methodology, and any tool that is important to that methodology, including the use of LLMs, should be described also. For example, students should mention tools (including LLMs) that were used for data processing or filtering, visualization, facilitating or running experiments, reviewing, coding, translating, proving theorems, etc. It may also be advisable to describe the use of LLMs in implementing the method (if this corresponds to an important, original, or non-standard component of the approach). Second, students are responsible for the entire content of their work, including all text and figures, so while students are welcome to use any tool they wish for writing, they must ensure that all text is correct and original.

A few practices are strongly discouraged: plagiarism, falsification, and any other methods that divert students' internal resources and abilities away from truly immersing themselves in solving this assignment.

So, the goal should not be to avoid challenges or merely achieve formal results, such as grades, but to engage deeply with the material and develop meaningful solutions.

Goals

The primary goals of this assignment are for students to:

1. Learn how to work with real-world data that may contain human labeling errors.
2. Understand how to design AI systems that account for uncertainty and doubt.
3. Explore methods and model configurations to enhance robustness when dealing with imperfect data.

To achieve these goals, students will compare three types of loss functions in the context of a classification problem: the *cross entropy loss* (CE loss), as a standard loss function for classification tasks; a *custom loss*, called B-loss (ideally, adapted or created by the student) that can estimate the uncertainty of the model's predictions; the *Heteroscedastic Regression Loss based on a multivariate normal distribution* (N-loss) which is normally used in fitting problems and which capable of estimating the uncertainty of the model's predictions in classification.

Through this comparison, students will gain insights into how different loss functions handle imperfect data and how to make models more resilient to such challenges.

Tasks and Requirements

The following points should be met:

1. Fill your name in the form above
2. Review the dataset CIFAR-10N and related work
3. Check the proposed methodology on loss functions, ensembling and metrics
4. Implement the proposed methodology in code for training an ensemble with a given loss and a list of seeds
5. Implement the label smoothing techniques and alternative ways to calculate final ensemble predictions using majority voting and weighted predictions (check your personal task with an instructor)

6. Perform a computational experiment, fill the table and figure of the results
7. Fill in the remaining sections of the report, including Related work, Results and discussion, Conclusion, and Introduction
8. Submit the code and report

Bonus (complete one of the following points as a bonus task):

1. Revise the proposed B-model or N-model, enhance it and obtain higher metrics in additional experiments
2. Implement a more advance ANN, e.g. using recommendations from papers related to CIFAR-10N dataset
3. Achieve the original results with accuracy of top 10 solutions in accordance with SOTA rating for CIFAR-10N

Sections that require revision and completion: 1, 2, 4, 5, 6.

Evaluation Criteria

The work will be evaluated based on the following criteria:

1. Basic Level (65–75%):

- The main tasks have been completed.
- Verifiable results have been obtained.
- A well-formatted written report and the main code have been submitted.

2. Intermediate Level (76–89%):

- The approach to solving the problem and the results obtained demonstrate originality.
- Additional efforts have been made to achieve better results.
- The work is neatly presented and shows a deep understanding of the topic.

3. Advanced Level (90–100%):

- One of the *bonus track* questions has been solved, or results exceeding basic expectations have been achieved.
- The work demonstrates exceptional quality, creativity, and independence.

1 Introduction

Nowadays data science field has a lot of benchmark datasets that have been carefully prepared for the purpose of model, architectures and method comparisons. However working in this paradigm is delusional since there is no perfect data in the real world. The model can be easily overfit on a synthetic data, which it becomes when labeled in a specific manner. One of the ways to deal with this problem - is to introduce the noise into the dataset. Either adding it to the input or to the target. Wei et al. discussed approaches in generating a dataset with the noise. It's important to compare the synthetic noise with human generated and approaches to deal with it. Modern DS challenges often come with the so-called "expert uncertainty problem", when the same level expert in the field have different opinions on how to label the data. For example the same x-ray image for one expert will contain a tumor of certain size and shape and for another expert no tumor at all. That is why it is important to create methodologies on how to deal with mislabeled data in order to improve the models accuracy. In this paper we address the issue of multiclass classification on CIFAR-10N dataset. Specifically, I compare the result of training a CNN model with 3 different kinds of loss-functions with label-smoothing technique and ensembled solution to analyze the best method to use.

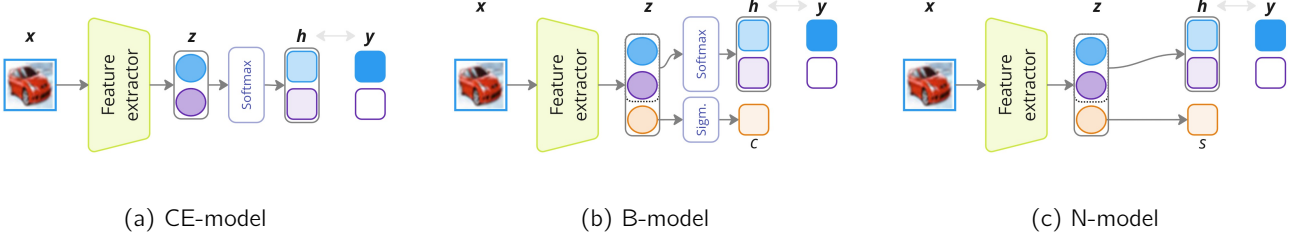


Figure 1: Schematics of the classification models under study: the CE-model serves as the baseline (see Equation (11)), the B-model includes an additional output $c \in (0, 1)$ that calculates the certainty value of the prediction \mathbf{h} (see Equation (7)), and the N-model introduces an additional output $s \in (-\infty, \infty)$ that computes the logarithmic variance of the prediction (see).

2 Related work

Learning with Real-World Human Annotations. Please analyze in more detail 2–4 papers related to the CIFAR-10N dataset [24], focusing on models and methods that can be implemented in this work for comparison or to enhance the results. Emphasis should be placed on approaches that have the potential to improve the outcomes of this study. Combine a descriptive approach with a critical analysis. Additionally, ensure that all methods, ideas, and concepts that are not your own are properly referenced.

Uncertainty-aware objectives. One of the approaches to uncertainty estimation of the regression models is the *heteroscedastic* regression that takes both the variable mean and variance into account [18, 20]. So, the model trains to predict means and variances, and the uncertainty of the model predictions can be estimated using the variance values. Fortunately, classification models can also use a *squared error* (SE) loss. Hui and Belkin [10] demonstrated that the SE and CE-based computer vision models are close in accuracy. However, a SE loss needs some more training epochs. Kendall and Gal [11] dealt with two types of uncertainty, that are aleatoric (data uncertainty) and epistemic (model uncertainty), and proposed two approaches in uncertainty estimation. Kendall and Gal [11] declared that out-of-data examples cannot be identified with aleatoric uncertainty. The authors also proposed an approach that combines aleatoric and epistemic uncertainties. Further work by van Amersfoort et al. [22] deals with the deterministic uncertainty quantification method. The proposed model learns the positions of centroids of classes and trains kernels to estimate the distance between an input sample and centroids, which allows the inference model to recognize an out-of-data sample as uncertain. Sensoy et al. [21] developed a theory of evidence perspective and represented the model predictions as a Dirichlet density distribution over the softmax outputs and proposed a novel loss function. Collier et al. [3] proposed a method for training deep classifiers under heteroscedastic label noise. The method deals with the softmax temperature tuning that allows to control a bias-variance trade-off.

Ensembling, test-time augmentation, and label smoothing. Ashukha et al. [1] demonstrated that many ensembling techniques are equivalent to an ensemble of several independently trained networks in terms of test performance. Test-time augmentation is a technique that improves model performance using averaging the predictions [14]. Probably the simplest ways to make models be more robust to noise in labels are label smoothing [23] and data augmentation [19].

Data uncertainty estimation in practice Corrupted inputs [13] and corrupted labels [26], in-domain and out-of-domain distributions [15, 11, 3] are some of the poles of the research in the scope of data uncertainty estimation. The typical test of the models in practice is to use public datasets with corrupted (noisy) labels at the training and validation stages but with clean labels at the test stage [25, 26]. A number of methods try to detect input samples with incorrect labels and remove [3, 25, 26] or under-weight these samples [11, 4]. Han et al. [8] declared that models learn data with clean labels first and noisy labels then, and proposed a new paradigm called *co-teaching* with the training of two networks.

3 Methodology

In general, consider a model $\mathbf{f}[\mathbf{x}, \mathbf{w}]$ parameterized by weights \mathbf{w} that maps an input \mathbf{x} into logits \mathbf{z} first and then into the hypothesis \mathbf{h} that approximates the ground truth \mathbf{y} . The negative log-likelihood minimization [18, 2, 6]

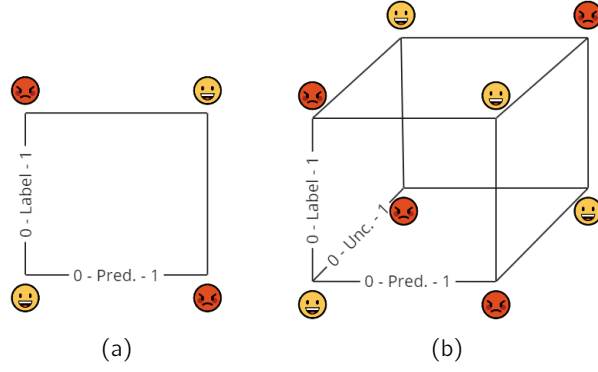


Figure 2: Binary classification intuition with the BCE loss (a) and the proposed binary B-loss eq. (4) (b) with respect to the values of the model's outputs (uncertainties $u = 1 - c$, predictions h), and labels y .

allows the formalization of the following uncertainty-aware loss functions for fitting and classification problems using different types of distributions for the outputs of the models.

3.1 B-loss

This paragraph presents a specific interpretation of a binary classification model based on minimizing the *uncertainty-aware negative log-likelihood with the Bernoulli distribution* (B-model, B-loss). The proposed model is trained to ensure that true predictions are made certain and false predictions if they occur, are made uncertain (see fig. 2). The binary classifier estimates the certainty value $c \in (0, 1)$, which is the primary task in the proposed formalization. Additionally, the classifier estimates and enhances the similarity δ between the hypothesis \mathbf{h} and the ground truth \mathbf{y} , which constitutes the secondary task in the proposed formalization.

Binary classification. Consider an i^{th} sample and a model with logits $z^{(i)} = [z_{pred}^{(i)}, z_{cert}^{(i)}]$ which correspond to a prediction $h^{(i)} = \sigma(z_{pred}^{(i)})$, and the certainty $c_i = \sigma(z_{cert}^{(i)})$ associated with the prediction, respectively. Then, compare the prediction $h^{(i)}$ and the given label $y^{(i)}$ using a scalar product metric $\delta_i = y^{(i)} h^{(i)}$, and map this metric as a pseudo-label of a binary uncertainty estimator into the parameters of a Bernoulli probability mass function [18]:

$$p_i = p(\delta_i | c_i) = \begin{cases} 1 - c_i & \text{if } \delta_i \rightarrow 0, \\ c_i & \text{if } \delta_i \rightarrow 1, \end{cases} \quad (1)$$

where $\delta_i \in (0, 1)$ is the smoothed pseudo-label that characterizes the similarity between the label and the prediction.

eq. (1) is a discrete probability distribution for a random variable that takes the value 0 with probability $1 - c_i$, which is an incorrect prediction that corresponds to an uncertainty of the prediction, and the value 1 with probability c_i , which is a right prediction that corresponds to a certainty of the prediction. The Bernoulli distribution has an equivalent power law form [18]:

$$p_i = c_i^{\delta_i} (1 - c_i)^{1 - \delta_i}. \quad (2)$$

For a roll-out of dataset of m i.i.d. pairs $\{x^{(i)}, y^{(i)}\}$ associated with the outputs of the model $\{h^{(i)}, c_i\}$, the joint probability [2] for the given probability mass function eq. (2) takes the following form:

$$P(\delta_1, \dots, \delta_m | c_1, \dots, c_m) = \prod_{i=1}^m c_i^{\delta_i} (1 - c_i)^{1 - \delta_i}. \quad (3)$$

The negative logarithm of the joint probability eq. (3) represents the proposed uncertainty-aware B-loss for the binary classification:

$$\mathcal{L}_B = -\frac{1}{m} \sum_{i=1}^m [\delta_i \log c_i + (1 - \delta_i) \log(1 - c_i)]. \quad (4)$$

eq. (4) intuition is demonstrated in fig. 2. The B-loss can be generalized for the case of multiclass classification.

Multiclass (N-classes) classification. Consider an i^{th} sample and a model with logits $\mathbf{z}^{(i)} = [\mathbf{z}_{pred}^{(i)}, z_{cert}^{(i)}]$ which correspond to a vector of prediction $\mathbf{h}^{(i)} = \text{softmax}(\mathbf{z}_{pred}^{(i)})$, $\mathbf{h}^{(i)} \in \mathcal{R}^N$, and the certainty $c_i = \sigma(z_{cert}^{(i)})$ associated with the prediction, relatively. Then, compare the prediction vector $\mathbf{h}^{(i)}$ and the given one-hot encoded label vector

$\mathbf{y}^{(i)}$ using a scalar product terms $\delta_k^{(i)} = y_k^{(i)} h_k^{(i)}$, and map this metrics as pseudo-labels into a probability mass function:

$$p_i = \prod_{k=1}^N \left(\frac{C_i}{N} \right)^{\delta_k^{(i)}} \left(\frac{1 - C_i}{N} \right)^{1 - \delta_k^{(i)}}, \quad (5)$$

where $\delta_k^{(i)} \in (0, 1)$ is the smoothed one-hot encoded pseudo-label that characterizes the similarity between the k^{th} components of the label and the prediction vectors, N is the number of classes.

Following the logical sequence given in section 3.1 and in [18], the joint probability for eq. (5) can be obtained, and then transformed into the negative log-likelihood (NLL):

$$NLL = -\frac{1}{m} \sum_{i=1}^m \left(\cos(\mathbf{h}^i, \mathbf{y}^i) \log \left(\frac{C^i}{N} \right) + (N - 1) (1 - \cos(\mathbf{h}^i, \mathbf{y}^i)) \log \left(\frac{1 - C^i}{N} \right) \right), \quad (6)$$

where $\cos(\mathbf{h}^i, \mathbf{y}^i)$ is the smoothed pseudo-label that characterizes the cosine similarity between two N -dimensional vectors: the vector of prediction and the one-hot encoded label vector.

Finally, the proposed uncertainty-aware B-loss for the N -classes classification is the Kulback-Loeberg divergence between two distributions: the one-hot encoded smoothed pseudo-labels distribution and the NLL distribution eq. (6):

$$\mathcal{L}_B = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^N \delta_k^{(i)} \log \delta_k^{(i)} + NLL. \quad (7)$$

where m is the number of samples (in a batch), N is the number of classes, $\delta_i = y^{(i)} h^{(i)}$ are the terms of a scalar product of the one-hot encoded label vector and the vector of prediction of the model, C^i is the certainty of the prediction (fig. 4b).

3.2 N-loss

Since the binary classification can be assumed as a particular case of a multiclass classification, this section skips the binary classification paragraph.

Multiclass classification. Consider an i -th sample and a model with logits $\mathbf{z}^{(i)} = [\mathbf{z}_{mean}^{(i)}, \mathbf{z}_{var}^{(i)}]$ which maps into the parameters of a multivariate normal distribution: the hypothesis or mean $\mathbf{h}^{(i)} = \mathbf{z}_{mean}^{(i)}$ that approximates the ground truth $\mathbf{y}^{(i)}$, and the variance $\sigma_{(i)}^2 = \exp(\mathbf{z}_{var}^{(i)})$ that characterizes the uncertainty of the hypothesis, $f[\mathbf{x}^{(i)}, \mathbf{w}] = [\mathbf{h}^{(i)}, \sigma_{(i)}^2]$. In other words, it is assumed that the conditional probability distribution $p = p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}) = p(\mathbf{y}^{(i)} | \mathbf{f}[\mathbf{x}^{(i)}, \mathbf{w}])$ has the form of a multivariate normal distribution characterized by equal variances (spherical covariances) in N -dimensional space [17]:

$$p^{(i)} = \frac{\exp \left(-\frac{\sum_{k=1}^N (y_k^{(i)} - h_k^{(i)})^2}{2\sigma_{(i)}^2} \right)}{(2\pi\sigma_{(i)}^2)^{\frac{N}{2}}}, \quad (8)$$

The multivariate normal distribution of (8) can be applied to the negative log-likelihood criterion of an uncertainty-aware negative log-likelihood loss (N-loss) for the regression [18]:

$$\mathcal{L}_N = \frac{1}{2m} \sum_{i=1}^m \left(\sum_{k=1}^N \frac{(y_k^{(i)} - h_k^{(i)})^2}{\sigma_{(i)}^2} + N(s^{(i)} + r) \right), \quad (9)$$

where m is the number of samples (in a batch), $\mathbf{y}^{(i)}$, $s^{(i)} = \log \sigma_{(i)}^2$ is the log-variance, $r = \log 2\pi$ is the constant value.

The last term in (9) represents a constant that can be neglected. Kendall et al. [11] recommended to train the models to predict log-variances $s^{(i)} = \log \sigma_{(i)}^2$, because it is more numerically stable than the variance $\sigma_{(i)}^2$ and the loss avoids a potential division by zero:

$$\mathcal{L}_N = \frac{1}{m} \sum_{i=1}^m \left(e^{-s^{(i)}} \sum_{k=1}^N (y_k^{(i)} - h_k^{(i)})^2 + Ns^{(i)} \right). \quad (10)$$

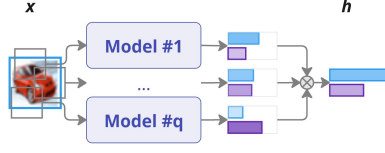


Figure 3: Ensemble of q models makes the prediction \mathbf{h} for the q augmented copies of an input sample \mathbf{x} . Two-classes classification is demonstrated.

Thus, (10) represents a heteroscedastic regression loss [18, 11], generalized for the case of a space of N dimensions. Our proposal is to use this loss for classification problems.

The baseline loss in a multi-class classification problem is the cross-entropy loss [6, 18, 2]:

$$L_{CE} = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^N y_k^{(i)} \log h_k^{(i)}. \quad (11)$$

3.3 Ensembling

A set of q models $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_q$ with different random initialization of the weights are trained with the same dataset. This aggregation reduces overfitting and provides more robust estimates by averaging out individual model errors [1]. Each j^{th} model predicts a class index for the given i^{th} input:

$$\hat{y}^{(i,j)} = \arg \max_k h_k^{(i,j)}, \quad (12)$$

where i, j, k are the dummy indexes that refer to j^{th} augmented version of i^{th} sample, and k^{th} component of the prediction vector or class index, $i \in (1, m), j \in (1, q), k \in (0, N - 1)$.

The final ensemble prediction class is typically determined by *majority voting* based on the class predictions of the individual j^{th} model.

$$\hat{y}^{(i)} = \text{mode}(\hat{y}^{(i,1)}, \hat{y}^{(i,2)}, \dots, \hat{y}^{(i,q)}). \quad (13)$$

The final ensemble prediction class can also be determined by *confidence-based weighted predictions* (see fig. 3). Each model predicts a class $\hat{y}^{(i,j)}$ and provides a confidence value $co^{(i,j)}$ for its prediction:

$$co^{(i,j)} = \max_k (h_k^{(i,j)}). \quad (14)$$

The aggregated confidence for class k is:

$$co_k^{(i)} = \sum_{j=1}^q co^{(i,j)} \cdot \mathbb{I}(\hat{y}^{(i,j)} = k), \quad (15)$$

where $\mathbb{I}(\cdot)$ is the indicator function, which returns 1 if the condition is true and 0 otherwise, $co_k^{(i)}$ is the total confidence for class k across all models.

The final predicted class is:

$$\hat{y}^{(i)} = \arg \max_k co_k^{(i)}. \quad (16)$$

The uncertainty estimation in deep ensembling is derived from the variance of the individual model predictions. Higher variance among the models' outputs indicates greater uncertainty, providing a measure of epistemic uncertainty.

3.4 Metrics

The standard classification metrics for the balanced datasets are the accuracy, receiver operating characteristic - area under curve (ROC-AUC) [2, 6]. A set of more specific metrics used in uncertainty quantification involves estimation the confidence (see eq. (14)) [15]: the Brier score, the entropy, the expected calibration error (ECE), the negative log-likelihood (NLL), the prediction interval coverage probability (PICP), the sharpness, etc. [16, 7, 5, 9].

Since the proposed B-model (see eq. (7)) and N-model (see eq. (10)) have an extra output, the following additional *certainty* metrics can be met:

Table 1: Accuracy (%) of models trained on the CIFAR-10N dataset with clean and noisy labels, as well as with label smoothing (LS). Each model was trained using seven different weight initialization seeds for 20 epochs, then ensembled. The mean test accuracy of individual models is compared with the accuracy of the ensemble using majority voting (EMV) and the ensemble with weighted predictions (EWP). The best results are highlighted in bold.

Method	Arch.	#Param. (train.par.)	LS	Accuracy (single model / EMV / EWP),%	
				Clean	Noisy (rand1)
Baseline CE loss	9-l.CNN	4.4 M (all)	0.4	87.23 ± 0.0035/90.14/00.00	81.18 ± 0.0071/84.70/00.00
Proposed B-loss				85.76 ± 0.0026/89.77/89.75	81.20 ± 0.0034/85.40/85.44
Proposed N-loss				85.88 ± 0.005/89.85/89.82	80.98 ± 0.0058/84.78/84.81

- $c^{(i)} \in (0, 1)$ for the B-model ;
- $1 - \text{sigm}(s^{(i)}) \in (0, 1)$ for the N-model.

Both of the above metrics can be used as weights in eq. (16), thus the *certainty-based weighted predictions* should be met.

4 Results and Discussion

CIFAR-10N dataset [24] was split into training, validation, and test sets in the amounts of [45000, 5000, 10000] samples, respectively. The models were trained with a 9-layer convolutional neural network [8, 25]. The network has 4.4 million parameters that were randomly initialized during training. The CNN architecture and most of the settings correspond to the experiments by Xia et al. [25] with minor changes: the models were trained for 20 epochs (200 in the original paper) using the Adam optimizer with a momentum of 0.9 and a batch size of 128, and a constant learning rate of 0.001 (in the original paper, the initial learning rate linearly decreased to zero starting from the 80th epoch); image samples were transformed into tensors and normalized with means of [0.491, 0.482, 0.447] and standard deviations of [0.247, 0.243, 0.261].

Some of the experimental settings might differ from those of Xia et al. [25]: the model ensembling technique was applied; a random resized crop with a scale range of [0.8, 0.1] and an aspect ratio range of [0.9, 1.1] was applied to all samples in all sets as a transformation, so test set sampling was implemented using test-time augmentation with a random resized crop; models with the lowest validation loss were used for inference.

All the experiments were performed seven times with random seeds of [42, 0, 17, 9, 3, 16, 2]. The mean and standard deviation of experimental results were then reported. The multiple predictions obtained through model ensembling allowed for the calculation of final predictions using majority voting.

Accuracy was used [12] as metric. The obtained results were not compared with the state-of-the-art results. The latter aggregate multiple techniques and complex network architectures. So, the comparison would be unfair.

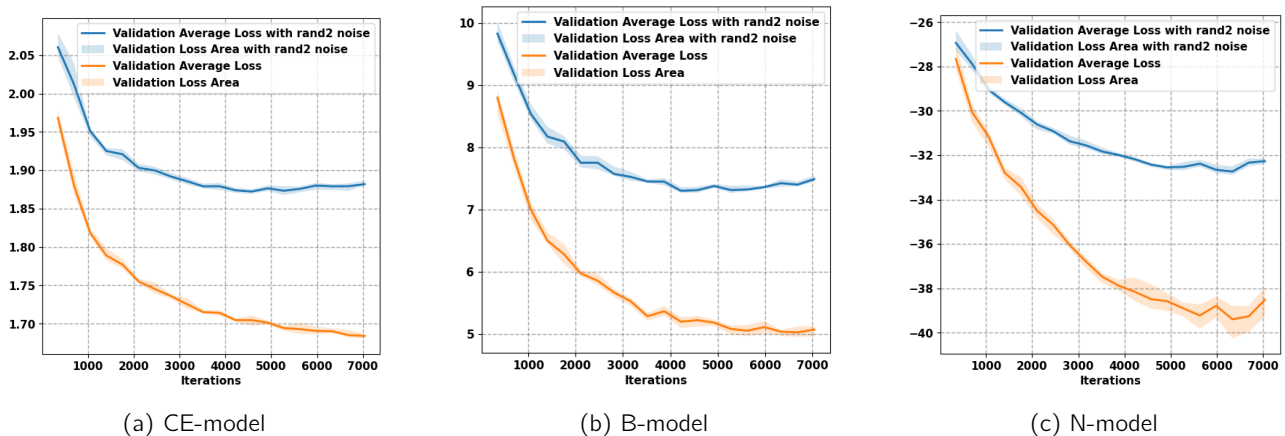


Figure 4: Validation loss values during the training of the ensemble models in 20 epochs with clean and noisy data: CE-model (a), B-model (b), and N-model (c).

Clean labels consistently yield higher accuracy scores, a result anticipated at the inception of the study. This is evident in the validation loss plots, where two distinct lines are clearly discernible: one representing the loss on clean labels and the other on noisy labels.

For the ensemble with weighted predictions, I utilized weights based on certainty. This method requires additional information—specifically, the model's confidence level, which can be derived from B-loss or calculated as $1 - \text{sigm}(s^{(i)}) \in (0, 1)$ for the N-model. Due to the limitations of cross-entropy (CE) loss, it was not feasible to implement weighted ensembles with this type of loss, and therefore, no results for CE-weighted ensembles are included in the table.

From the data presented in the table, B-loss consistently produced results with smaller deviation compared to other losses. Notably, weighted ensembles performed worse on clean labels relative to majority voting but demonstrated superior performance on noisy data. Furthermore, ensemble methods generally outperformed single models, regardless of the type of ensemble used. This finding underscores the effectiveness of ensembles in enhancing accuracy.

Overall, cross-entropy loss proved to be the optimal solution for clean labels, with other losses yielding slightly inferior results. In contrast, for the noisy dataset, the differences among the losses were less pronounced, but B-loss emerged as the leading approach. These findings suggest that ensemble methods, particularly those incorporating weighted predictions based on model certainty, are a robust strategy for improving accuracy in both clean and noisy data environments.

5 Ethics

For this study, several large language models (LLMs) were utilized during the experimental phase to assist in debugging errors. Specifically, we used LLama 3.2 and DeepSeek R1 models. Additionally, as a final refinement, the Qwen 2.5 model was employed to enhance the text quality. Unless explicitly stated by the authors, no further contributions from these models were made towards the generation of ideas, writing the code, or similar tasks.

This project was conducted in the context of a course homework assignment, involving a group of students who collectively used approximately $12 \cdot N_{students}$ hours of computational time. This significant computational effort has generated a carbon footprint that should be considered by the ethical committee when reviewing

6 Conclusion

Modern real-world problems often come with inherent uncertainties, necessitating robust techniques to address these challenges. This study demonstrates that ensembles remain one of the most effective methods for improving model accuracy. Additionally, several approaches can be employed to introduce uncertainty into models, which not only slightly enhances accuracy but also serves as a valuable tool for decision-making.

When working with human-labeled data, the uncertainty in the model can highlight specific cases that require additional expert labeling, thereby improving dataset quality.

For future work, I suggest experimenting with an unlabeled dataset from a computer vision (CV) problem where domain experts have differing opinions. The goal is to design a system for trial and relabeling of images where the model exhibits high uncertainty, aiming to show how this process can converge to stable and highly accurate models. This hypothesis remains unproven, and the design of losses and the overall relabeling process will significantly influence the results. Therefore, these aspects must be carefully analyzed and documented.

References

- [1] A. Ashukha, A. Lyzhov, D. Molchanov, and D. Vetrov. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. *arXiv preprint arXiv:2002.06470*, 2020.
- [2] C. M. Bishop and N. M. Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [3] M. Collier, B. Mustafa, E. Kokiopoulou, R. Jenatton, and J. Berent. A simple probabilistic method for deep classification under input-dependent label noise. *arXiv preprint arXiv:2003.06778*, 2020.
- [4] E. Engleson, A. Mehrpanah, and H. Azizpour. Logistic-normal likelihoods for heteroscedastic label noise, 2023.
- [5] T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [6] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.

- [7] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. *International Conference on Machine Learning*, pages 1321–1330, 2017.
- [8] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels, 2018.
- [9] F. Hernandez, L. Bertino, G. Brassington, E. Chassignet, J. Cummings, F. Davidson, M. Drevillon, G. Garric, M. Kamachi, J. M. Lellouche, et al. Probabilistic forecasting in meteorology: A review. *Quarterly Journal of the Royal Meteorological Society*, 141(688):318–350, 2015.
- [10] L. Hui and M. Belkin. Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks, 2021.
- [11] A. Kendall and Y. Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- [12] A. Kumar, P. Liang, and T. Ma. Verified uncertainty calibration, 2020.
- [13] E. Mintun, A. Kirillov, and S. Xie. On interaction between augmentations and corruptions in natural corruption robustness, 2021.
- [14] D. Molchanov, A. Lyzhov, Y. Molchanova, A. Ashukha, and D. Vetrov. Greedy policy search: A simple baseline for learnable test-time augmentation, 2020.
- [15] T. Pearce, A. Brintrup, and J. Zhu. Understanding softmax confidence and uncertainty. *CoRR*, abs/2106.04972, 2021. URL <https://arxiv.org/abs/2106.04972>.
- [16] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [17] S. Prince. *Computer Vision: Models Learning and Inference*. Cambridge University Press, 2012.
- [18] S. J. Prince. *Understanding Deep Learning*. MIT Press, 2023. URL <http://udlbook.com>.
- [19] S.-A. Rebuffi, S. Goyal, D. A. Calian, F. Stimberg, O. Wiles, and T. A. Mann. Data augmentation can improve robustness. *Advances in Neural Information Processing Systems*, 34:29935–29948, 2021.
- [20] M. Seitzer, A. Tavakoli, D. Antic, and G. Martius. On the pitfalls of heteroscedastic uncertainty estimation with probabilistic neural networks. *arXiv preprint arXiv:2203.09168*, 2022.
- [21] M. Sensoy, L. Kaplan, and M. Kandemir. Evidential deep learning to quantify classification uncertainty, 2018.
- [22] J. van Amersfoort, L. Smith, Y. W. Teh, and Y. Gal. Simple and scalable epistemic uncertainty estimation using a single deep deterministic neural network. *CoRR*, abs/2003.02037, 2020. URL <https://arxiv.org/abs/2003.02037>.
- [23] J. Wei, H. Liu, T. Liu, G. Niu, M. Sugiyama, and Y. Liu. To smooth or not? when label smoothing meets noisy labels. *arXiv preprint arXiv:2106.04149*, 2021.
- [24] J. Wei, Z. Zhu, H. Cheng, T. Liu, G. Niu, and Y. Liu. Learning with noisy labels revisited: A study using real-world human annotations, 2022. URL <https://arxiv.org/abs/2110.12088>.
- [25] X. Xia, T. Liu, B. Han, M. Gong, J. Yu, G. Niu, and M. Sugiyama. Sample selection with uncertainty of losses for learning with noisy labels, 2021.
- [26] Q. Yao, H. Yang, B. Han, G. Niu, and J. Kwok. Searching to exploit memorization effect in learning from corrupted labels, 2020.