
Uncertainty Estimation Using a Single Deep Deterministic Neural Network - ML Reproducibility Challenge 2020

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 The investigated paper claims RBF network when trained with BCE loss along with two-sided gradient penalty
4 outperforms deep ensemble in the task of out of distribution(OoD) detection along with competitive accuracy to softmax
5 based models. The Paper claims to outperform AUROC on Fashion-MNIST vs MNIST and CIFAR-10 vs SVHN. The
6 proposed algorithm is reported to detect both aleatoric and epistemic uncertainty as OoD. Authors mention the need for
7 a formal way to distinguish between the two kinds of uncertainty and pose it as an interesting future research avenue.

8

9 The scope of this report is to reproduce the results related to OoD detection presented in the paper. Along
10 with the reproduction of results, we propose an extension for explicit detection of aleatoric and epistemic uncertainty as
11 intended by the authors.

12 **Methodology**

13 The author's training code is available on GitHub. Additionally, we have made available all the experimentation codes
14 in the form of notebooks. We provide all the results and analysis on the models described in the paper and trained on
15 NVIDIA Tesla T4 GPU.

16 **Results**

17 Overall our reproduction supports the claims of the paper, we can replicate trends and plots as described in the paper.
18 Most of the results are within 1% of the value reported. Notably, AUROC(M) of DUQ in OoD detection of Fashion-
19 MNIST vs MNIST is off by 1.5% and we got a different optimal value for gradient penalty weight (λ) in it.
20 Also, the results of our proposed extension and its analysis are encouraging. Our proposed extension provided an
21 increase of 1.8% in AUROC(M) in Fashion-MNIST vs MNIST and provided explicit control over the aleatoric and
22 epistemic uncertainty.

23 **What was easy**

24 The proposed approach is quite simple and elegant. The availability of the author's code made the implementation of
25 various experiments easy.

26 **What was difficult**

27 Understanding of proposed approach requires advanced knowledge of calculus related to the Lipschitz constant and its
28 role to quantify the upper bound of the sensitivity of any function.

29 **Communication with original authors**

30 We discussed our report with the authors, they find our analysis on aleatoric uncertainty interesting and appreciate our
31 proposed extension and its encouraging results.

32 1 Introduction

33 Reliable uncertainty quantification has been a challenge in deep learning, models giving overconfident wrong
34 predictions can not be deployed for practical purposes and if deployed can even be fatal. The investigated paper
35 Uncertainty Estimation Using a Single Deep Deterministic Neural Network [1] presents a method of training deep
36 neural networks to detect out of distribution(OoD) points in a single forward pass along with classification.
37

38 Authors define a set of feature vectors called centroids each corresponding to a certain class. Distance between feature
39 vectors predicted by model and centroids are used for class prediction. This architecture along with BCE loss employs
40 a two-sided-gradient penalty to make the model sensitive to input and hence providing OoD detection ability, this
41 approach is called deterministic uncertainty quantification, DUQ.

42 2 Scope of reproducibility

43 The proposed algorithm outperforms the state-of-the-art techniques for OoD detection on Fashion-MNIST vs MNIST,
44 CIFAR-10 vs SVHN in terms of AUROC as well as computational time. DUQ, an RBF network-based algorithm
45 provides competitive accuracy to softmax models along with strong OoD detection ability. The scope of this report is as
46 follows.

47 2.1 Addressed claims from the original paper

- 48 • OoD detection ability on various datasets as reported in paper.
- 49 • Role of different hyperparameters (Ablation Study).

50 2.2 Other experiments

- 51 • Noise sensitivity of proposed algorithm.
- 52 • Propose extension E-DUQ for explicit detection of aleatoric and epistemic uncertainty .

53 3 Methodology

54 3.1 Model descriptions

55 Figure 1 represents the complete algorithm, first, the features are computed through a standard feature extractor (f_θ)
56 after this extracted features are passed through class-specific layers (W_c) to calculate the feature vector for each class.
57 Distance (K_c) of this vector from the centroid in a kernel space represents uncertainty in prediction. The centroid is
58 calculated for each batch as described in the equations below.

$$59 n_{c,t} = \gamma \times n_{c,t-1} + (1-\gamma) \times n_{c,t}$$

60

$$m_{c,t} = \gamma \times m_{c,t-1} + (1-\gamma) \times \sum W_c f_\theta(x_{c,t,i})$$

$$61 e_{c,t} = m_{c,t} / n_{c,t}$$

62

63 Where $n_{c,t}$ is number of images of class c in a particular batch at time t , $m_{c,t}$ is weighted sum of feature vectors
64 through different mini batches and γ is the momentum for updating centroids. $x_{c,t,i}$ is the element i of minibatch at
65 time t of class c .
66

67 Total cost comprises of two-loss functions, one is the BCE loss which is the sum of cross-entropy of a binary one-hot
68 vector with the actual label of that class as described in the equation below.
69

$$70 L_1(x, y) = \sum_c y_c \log(K_c) + (1 - y_c) \log(1 - K_c)$$

71

72 Usually, deep learning models suffer from *feature collapse* where the presence of some non-robust discriminative
73 features in input space will map inputs to same outputs, the degree to which a model prevents this feature collapse is its

74 *sensitivity*, a two-sided gradient penalty loss is used to enforce sensitivity to our model which is defined as:

75

76 Two-sided GP loss = $\lambda \cdot \left[\|\nabla_x \sum_c K_c\|_2^2 - 1 \right]^2$

77

78 Which essentially keeps the norm of Jacobian above a threshold and prevents feature extractor from collaps-
79 ing to a constant function.

80 Another version of gradient penalty can be a single-sided loss as:

81

82 Single-sided GP loss = $\lambda \cdot \max(0, \|\nabla_x \sum_c K_c\|_F^2 - 1)$

83

84 The investigated paper claims two-sided gradient penalty can enforce sensitivity in the model while a single-sided
85 penalty is not able to do so. We will verify these claims empirically in sections to go.

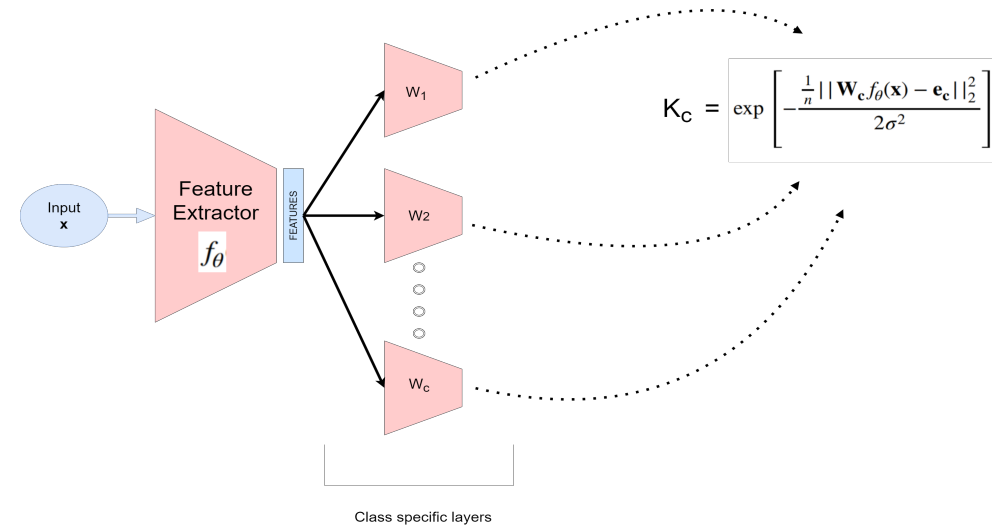


Figure 1: Figure describing complete algorithm, for each class it computes RBF score (K_c) which is the measure of OoD detection

86 **3.2 Datasets**

87 **3.2.1 Two-moons dataset**

88 It is a two-dimensional dataset. As the name suggests, it has two intertwined classes on the shape of the crescent. We
89 obtain it using Sciket-learn implementation as described by the authors.

90 **3.2.2 Fashion-MNIST and MNIST**

91 Fashion-MNIST is a dataset of 28x28 grayscale images consisting of a training set of 60,000 images and a test set of
92 10,000 images each belonging to either of 10 classes. MNIST is a database of 28x28 grayscale handwritten digit images
93 having a training set of 60,000 and a test set of 10,000 images. We obtain both datasets using Pytorch's dataset module.
94 As we are testing the proposed algorithm on out-of-distribution detection Fashion-MNIST, MNIST has been a notably
95 difficult pair for this task [2].

96 **3.2.3 CIFAR-10 and SVHN**

97 CIFAR-10 is a dataset of 60,000 32x32 colour images in 10 classes, having 50,000 training images and 10,000 test
98 images. SVHN is a real-world image dataset similar in flavor to MNIST but incorporates an order of magnitude more
99 labeled data with 604,388 training set images and 26,032 images for testing.
100 Similarly we will evaluate algorithm on CIFAR-10 vs SVHN OoD detection [2].

101 **3.3 Experimental setup**

102 The author’s training code is available on GitHub which was quite useful for us. We have added
103 our codes for all analysis and experimentation on [https://drive.google.com/drive/folders/](https://drive.google.com/drive/folders/1bBrn2jRnRTIhxrAi7BXL0PkzCuUKlizF?usp=sharing)
104 [1bBrn2jRnRTIhxrAi7BXL0PkzCuUKlizF?usp=sharing](https://drive.google.com/drive/folders/1bBrn2jRnRTIhxrAi7BXL0PkzCuUKlizF?usp=sharing). We trained all models on NVIDIA Tesla T4 GPU.
105

106 Execution of the code is not computationally extensive. In the described experimental settings training on MNIST takes
107 15 minutes while training on CIFAR-10 takes 4 hours. Maximum 3 GB of GPU and 4 GB of RAM is sufficient to run
108 the code.

109 We followed the Appendix of the reported paper for architecture, optimization, the number of epochs, batch size and
110 generating data in the case of two-moons.

111 Additionally, we also find and recommend applying early stopping in training on CIFAR10 at around 60 epochs to
112 avoid over-training.

113 **4 Results**

114 In this section, we provide our empirical analysis on claims in Section 2, first on toy dataset, two-moons, and then
115 extend the analysis for Fashion-MNIST and CIFAR-10. Our experiments verify the claims of paper with similar results
116 on these datasets.

117 **4.1 Two-moon**

118 Our implementation on two-moons for the deep ensemble, DUQ without gradient penalty, with single and with a
119 two-sided penalty are shown in Figs 2-5. The yellow colour represents the region of certainty while the blue region
120 shows uncertainty.

121 It can be seen DUQ with a two-sided penalty is certain only near the training data and uncertain elsewhere, hence
122 showing the best uncertainty estimation over the region as compared to other methods. These figures validate the
123 importance of the two-sided penalty in enforcing sensitivity.

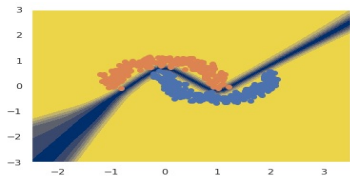


Figure 2: Uncertainty estimation using Deep Ensemble

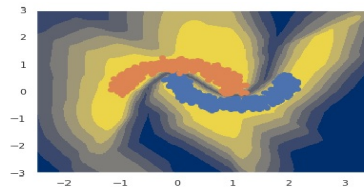


Figure 3: Uncertainty estimation using DUQ with $\lambda=0$

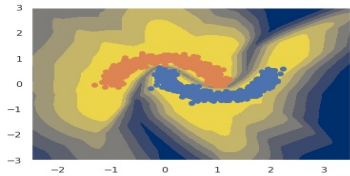


Figure 4: Uncertainty estimation using DUQ with single sided gradient penalty

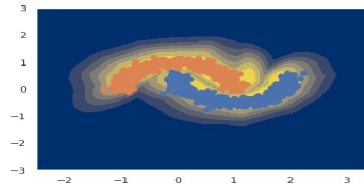


Figure 5: Uncertainty estimation using DUQ with two-sided penalty

124 **4.2 Fashion-MNIST**

125 DUQ is trained on Fashion-MNIST and its ability to detect MNIST datapoints as OoD is measured by AUROC, larger
126 being the better. In this section, we provide our results and analysis for DUQ on Fashion MNIST claimed by the
127 authors.
128

Method	Accuracy(FM)	AUROC(M)	Train time(in sec)	Inference time
Deep Ensemble	93.30 \pm 0.32	0.889 \pm 0.005	45	2.3
DUQ ($\lambda = 0.05$)	92.13 \pm 0.29	0.947 \pm 0.005	23	1
E-DUQ	92.35 \pm 0.15	0.964 \pm 0.005	23	1

Table 1: Result comparison of models trained on Fashion MNIST by different methods(mean over 3 runs), AUROC(M) is for separating Fashion-MNIST from MNIST. EDUQ is described in section 4.2.2.

129 In Table 1, we compare DUQ with deep ensemble, DUQ outperforms deep ensemble in AUROC as well as on training
130 and inference time. We observe a difference of 1.5% in AUROC of DUQ from the original paper.
131 In an algorithm that can detect OoD points, another measure of performance could be rejection classification, in which
132 we mix two-datasets and reject points based on uncertainty predicted by the model and expect an increase in accuracy.
133 In Fig 8 we find that rejection classification performance on Fashion-MNIST vs MNIST for DUQ in comparison with
134 deep ensemble gives more accuracy increase per rejection as the area below DUQ curve is more that area below DE.
135 These experiments broadly support the claims made by the authors.
136

137 4.2.1 Aleatoric sensitivity

138 Authors provide OoD detection analysis on Fashion-MNIST vs MNIST datapoints which are fundamentally different
139 datasets. The paper mentions that DUQ captures both aleatoric and epistemic uncertainty but doesn't provide any
140 analysis on these datasets.
141

142 To evaluate the noise detection sensitivity of DUQ, we plot the rejection classification performance of DUQ on
143 Fashion-MNIST vs Noisy Fashion-MNIST as shown in Fig 6. Noisy Fashion-MNIST is Fashion MNIST added with a
144 zero-centered gaussian noise with an std of 0.05 which is imperceptible to the human eye. We have used full test-dataset
145 (1:1 proportion) for the mix.
146 It can be seen that DUQ can even detect Noisy Fashion-MNIST datapoints with this low noise whereas deep ensemble
147 is unable to do so. This clearly states that DUQ is very sensitive to noise and deems even very low noise data
148 points as OoD which can even be undesirable for practical purposes, as extreme sensitivity will harm the model's
149 decision-making ability. This observation is due to the imposed sensitivity to the model by gradient penalty loss. We
150 address this issue in the next subsection and propose an extension.
151

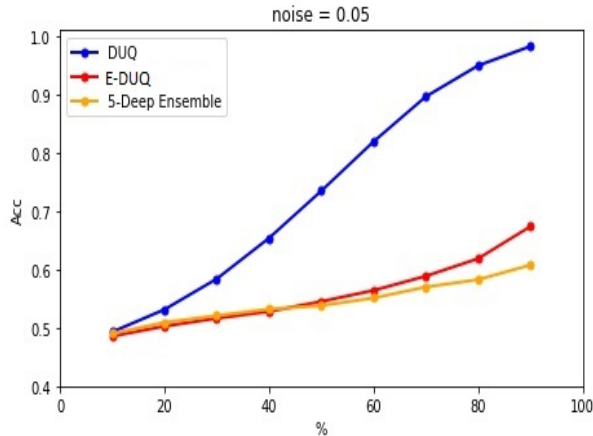


Figure 6: Rejection classification performance by different methods for Fashion-MNIST vs Noisy Fashion-MNIST. (X-axis is the % of data rejected using uncertainty estimates, Y-axis is the corresponding accuracy)

152 4.2.2 Learnable σ : Modelling noise through input

153 We propose to incorporate noise computation by predicting length scale(σ) for each data point rather than keeping it
154 constant as done by [3]. Authors of [3] showed how learning length scale can make model learn to attenuate loss from

155 erroneous labels. We name this extended approach Explicit-DUQ or E-DUQ as it gives us explicit control over aleatoric
 156 and epistemic uncertainty.

157

158 If noisy points are predicted with high σ , it will increase RBF score as shown in Fig 1 hence, disabling the model to
 159 predict noisy points as OoD. Hence E-DUQ should label only fundamentally different data distributions (i.e. epistemic
 160 uncertainty) as OoD. In E-DUQ, the feature extractor outputs features along with one extra logit for length scale, it will
 161 be very similar to the model described in Fig 1, but now along with $f_{\theta}(x)$ length scale(σ) will also be provided by the
 162 model to compute K_c .

163

164 Fig 7 shows the histogram of the magnitude of σ predicted by E-DUQ on Fashion-MNIST and Noisy Fashion-MNIST
 165 test set, it can be seen that model predicts higher values of σ for noisy data points.

166 Also, Rejection classification plot of E-DUQ (Aleatoric section, Fig 6) lies well below than that of DUQ which validates
 167 E-DUQ does not label noisy points as OoD and is robust to noise.

168

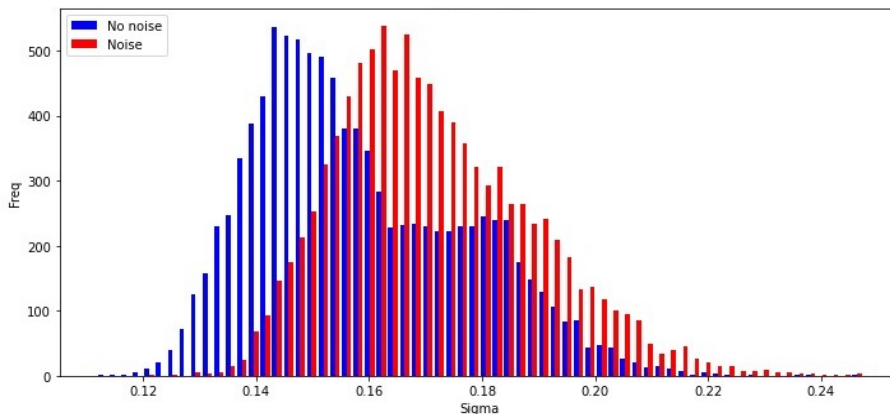


Figure 7: Histogram of values of length scale(sigma) predicted by E-DUQ for Fashion-MNIST and Noisy Fashion-MNIST

169 Table 1 illustrates E-DUQ outperforms DUQ trained on Fashion-MNIST in AUROC(M), this is because MNIST
 170 datapoints are OoD with respect to Fashion-MNIST in terms of epistemic uncertainty, which E-DUQ models explicitly.
 171 With the above empirical analysis, we can say it is possible to predict epistemic and aleatoric uncertainties independently
 172 using RBF score(K_c) and length scale(σ) respectively.

173 4.2.3 Other Hyperparameters and Ablations

174 In the following tables, ablations are done over centroid dimension, gradient penalty factor, and gradient penalty
 175 constant respectively (averaged over 5 runs). All other hyper-parameters are same as given in the original paper.

176 Table 2 shows how the performance of the model varies with the dimension of the centroid vector (z), we note AUROC
 177 first increases with z and then saturates while accuracy remains unaffected. Authors used $z = 256$ which according to
 178 us is a descent choice.

179

180 In Table 3, we show how the performance of DUQ varies with the value of λ , we observe first increase and then
 181 decrease in AUROC score which is also claimed by authors. We find a 1.5% deviation in the best AUROC(M) score as
 182 reported by the authors. They also mentioned best AUROC(NM) coincides with best AUROC(M) but we are not able to
 183 see this in our reproduction.

184 We also experiment how varying the constant number(α) in two-sided gradient penalty $\left[\|\Delta_x \sum_c K_c\|_2^2 - \alpha \right]^2$ affect
 185 training in Table 4. We fix $\lambda=1$ and experiment for different values for α , we find its performance to be better than
 186 model trained with optimal λ (in table 3). Thus, we infer that α can be a crucial hyperparameter.

187

188 Overall ablation study supports the claims of the paper with minor deviations, also we show some important hyperpa-
 189 rameters which are not stressed by the authors.

Size	Accuracy	AUROC(M)
10	92.15 \pm 0.15	0.887 \pm 0.015
256	92.13 \pm 0.29	0.947 \pm 0.005
500	92.32 \pm 0.15	0.959 \pm 0.004
1000	92.11 \pm 0.32	0.955 \pm 0.002

Table 2: Accuracy and AUROC(M) trend with centroid-vector dimension on Fashion MNIST.

λ	Accuracy	AUROC(NM)	AUROC(M)
0	92.39 \pm 0.08	0.941 \pm 0.005	0.941 \pm 0.011
0.05	92.13 \pm 0.29	0.962 \pm 0.007	0.947 \pm 0.005
0.1	92.11 \pm 0.08	0.953 \pm 0.006	0.942 \pm 0.007
0.2	92.03 \pm 0.11	0.929 \pm 0.006	0.954 \pm 0.005
0.3	92.11 \pm 0.22	0.941 \pm 0.016	0.946 \pm 0.008
0.5	91.94 \pm 0.16	0.939 \pm 0.008	0.931 \pm 0.014
1.0	91.12 \pm 0.23	0.895 \pm 0.036	0.901 \pm 0.029

Table 3: Accuracy and AUROC trend with λ for model trained on Fashion-MNIST.

α	Accuracy	AUROC(M)
0.01	91.89 \pm 0.03	0.928 \pm 0.002
0.1	92.23 \pm 0.08	0.951 \pm 0.004
0.2	92.36 \pm 0.13	0.953 \pm 0.004
0.5	91.15 \pm 0.19	0.918 \pm 0.005

Table 4: Accuracy and AUROC(M) trend with α in GP Loss on Fashion-MNIST

190 4.3 CIFAR-10

191 Till now, we have validated the central claim of the paper on the Fashion-MNIST dataset, in this subsection we provide
 192 few experiments to validate the claims made by authors about DUQ on CIFAR-10.

Method	Accuracy	AUROC	Train time(in sec)	Inference time(in sec)
Deep Ensemble	94.44 \pm 0.42	0.949 \pm 0.003	300	14
DUQ ($\lambda = 0.5$)	93.45 \pm 0.32	0.931 \pm 0.003	210	4

193 Result comparison on model trained on CIFAR-10 by deep ensemble and DUQ, AUROC reported is for separating
 194 CIFAR-10 from SVHN (averaged over 3 runs)

195 In Table 4.3, we compare DUQ with deep ensemble trained on CIFAR-10, the deep ensemble performs slightly
 196 better(1%) than DUQ in terms of accuracy and AUROC however DUQ is better in terms of computational time. We get
 197 similar values as in the original paper with a deviation of 0.5% in the values.

198 We obtain a similar rejection classification plot for CIFAR-10 vs SVHN as reported in the paper i.e. we see both DUQ
 199 and deep ensemble capture equal area as shown in Fig 9.

200

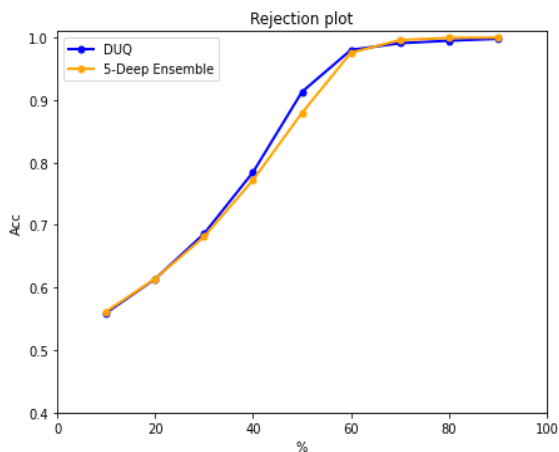


Figure 8: Rejection classification plot for deep ensemble and DUQ for Fashion-MNIST vs MNIST, x-axis is the % of data rejected using uncertainty estimates. Y-axis is the corresponding accuracy of prediction.

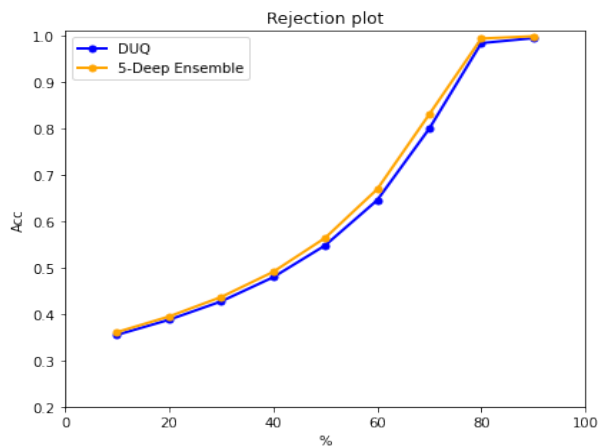


Figure 9: Rejection classification plot for deep ensemble and DUQ on CIFAR-10 vs SVHN. (X-axis is the % of data rejected using uncertainty estimates. Y-axis is the corresponding accuracy)

201 We provide a prediction histogram of DUQ on CIFAR-10 vs SVHN as shown in Fig 10. As datasets are of different
 202 sizes, we made the frequencies of histograms in proportion to dataset size as done by authors. We get the same graph as
 203 in the paper where the plot for CIFAR-10 is skewed towards the right extreme depicting higher certainty whereas a
 204 uniform plot for SVHN indicates uncertainty. These results support the claims on CIFAR-10 from the paper.

205

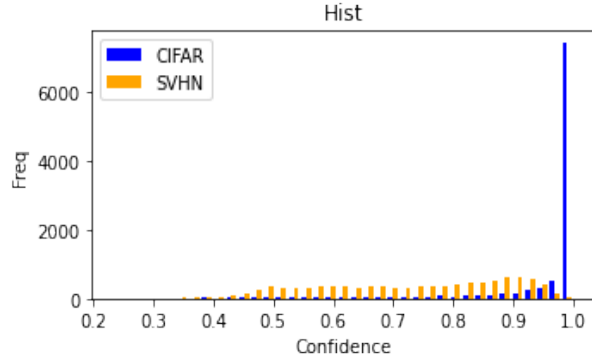


Figure 10: Uncertainty histogram of CIFAR-10 and SVHN for DUQ ($\lambda = 0.5$) trained on CIFAR-10

206 Additionally, we provide an analysis of the behavior of DUQ on noisy CIFAR-10 data points. In Figure 11, Noisy
 207 CIFAR-10 is CIFAR-10 added with a zero-centered gaussian noise with an std of 0.1. And all plots are made in propor-
 208 tion to the dataset size. We run the experiment with both the models and can see that DUQ is unconfident with Noisy
 209 CIFAR-10 points whereas deep ensemble can classify them with certainty, this finding is in line with that of Section 4.2.1.
 210

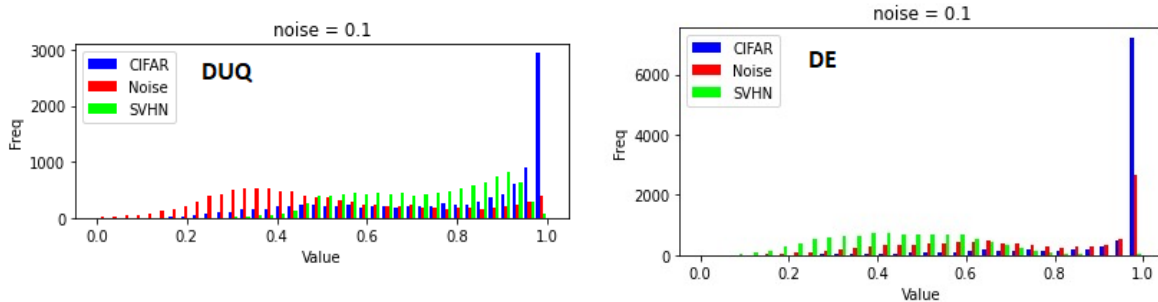


Figure 11: Uncertainty histogram of CIFAR-10 , SVHN and Noisy CIFAR-10 for DUQ ($\lambda = 0.5$) and deep-ensemble trained on CIFAR-10

211 5 Discussion

212 A simple approach along with data and models involved being in our computational limits helped us in the reproduction
 213 of results. We can validate most of the results and trends as reported in the paper.

214 Section 4.2 shows the overall behavior of DUQ and validates most of the claims made by authors like the importance
 215 of two-sided gradient penalty in the loss for enforcing the desired sensitivity which can show the state of the art
 216 OoD detection performance with a modest RBF network. Our additional experiments for understanding the nature of
 217 sensitivity of DUQ show how the addition of noise which is imperceptible to human eyes makes DUQ vulnerable for
 218 practical purposes.

219 Paper mentioned work on aleatoric and epistemic uncertainty detection is required, our naive extension of predicting σ
 220 inspired by [3] can even outperform the author’s results on Fashion-MNIST along with explicit control over uncertainty
 221 detection which is required for practical purposes. These preliminary results are encouraging for future research to be
 222 done in this direction.

223 Section 4.3 shows the scalability of different analyses for DUQ on larger datasets.

224 Overall we conclude that our reproduction results support different claims by authors of DUQ, it is an easy to implement,
 225 time-efficient algorithm with high sensitivity to input making it a good choice for OoD detection.

226 **References**

- 227 [1] Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep
228 deterministic neural network. 2020.
- 229 [2] Eric T. Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Görür, and Balaji Lakshminarayanan. Hybrid models
230 with deep and invertible features. In *ICML*, pages 4723–4732, 2019.
- 231 [3] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In
232 I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in*
233 *Neural Information Processing Systems*, volume 30, pages 5574–5584. Curran Associates, Inc., 2017.