
Reproducing Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning

Lana Baric
Eurecom
lana.baric@eurecom.fr

Abstract

This project reproduces the key findings of Gal and Ghahramani (2016), who proposed a Bayesian interpretation of dropout in deep neural networks. Specifically, they showed that applying dropout at both training and test time enables approximate variational inference and allows for modeling of epistemic uncertainty. To evaluate this claim, I implemented Monte Carlo Dropout (MC Dropout) in a variety of settings, including regression, classification, reinforcement learning, and predictive performance benchmarking. My experiments confirmed that MC Dropout produces meaningful uncertainty estimates without architectural changes or high computational cost. Despite certain constraints—such as limited dataset access and reduced computation budgets, the results align with the original paper’s conclusions and support MC Dropout as a scalable and practical tool for uncertainty-aware deep learning.

1 Introduction

Modern deep learning models have achieved impressive performance across tasks such as image recognition, natural language processing, and control. Despite this success, most deep neural networks (DNNs) provide point predictions without reliable estimates of uncertainty, which is problematic in safety-critical or data-scarce applications. In contrast, Bayesian models offer a principled framework for capturing uncertainty by maintaining distributions over model parameters, but they often come with high computational cost and complex inference.

In their work, Gal & Ghahramani [1] proposed a reinterpretation of dropout, a popular regularization technique in DNNs, as a form of approximate variational inference in a Bayesian neural network. They showed that by applying dropout at both training and test time, one can obtain Monte Carlo (MC) estimates of predictive uncertainty — enabling deep models to express confidence in their outputs without architectural changes or expensive sampling techniques.

My task was to try and replicate their study. I implemented and evaluated key experiments from their work, including: Model uncertainty in regression and classification, model uncertainty in reinforcement learning, and predictive performance. The goal was to better understand the behavior of MC Dropout in more practical settings and see how well dropout-trained networks approximate Bayesian inference and quantify their ability to capture meaningful predictive uncertainty.

2 Theoretical Background

Dropout is widely used in neural networks as a regularization method that prevents overfitting by randomly deactivating units during training. The authors of the original paper provided a probabilistic interpretation of dropout, showing that it approximates variational inference in a Bayesian deep Gaussian process (GP) model. This reinterpretation enables the use of dropout-trained networks to estimate predictive uncertainty in both regression and classification tasks.

2.1 Dropout as Variational Inference

Let $\omega = \{W_1, \dots, W_L\}$ be the set of weights in an L -layer neural network. The standard dropout objective includes a data-fitting term and L2 regularization:

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

In the Bayesian interpretation, dropout defines a variational distribution $q(\omega)$ where each weight matrix is masked by a random binary variable:

$$W_i = M_i \cdot \text{diag}(z_i), \quad z_{ij} \sim \text{Bernoulli}(p_i) \quad (2)$$

Here, M_i are trainable parameters and z_i are sampled dropout masks. This induces a highly multi-modal posterior distribution over weights, and training with dropout minimizes the KL divergence between $q(\omega)$ and the true posterior $p(\omega \mid \mathcal{D})$:

$$KL(q(\omega) \parallel p(\omega \mid \mathcal{D})) \quad (3)$$

2.2 Predictive Uncertainty with MC Dropout

At inference time, instead of turning dropout off, we keep it active and perform T stochastic forward passes to sample from $q(\omega)$. This Monte Carlo approximation yields predictive estimates:

$$\begin{aligned} \mathbb{E}[y^*] &\approx \frac{1}{T} \sum_{t=1}^T \hat{y}_t^* \\ \text{Var}[y^*] &\approx \frac{1}{T} \sum_{t=1}^T \hat{y}_t^* \hat{y}_t^{*T} - \mathbb{E}[y^*]^2 + \tau^{-1} I \end{aligned}$$

The precision parameter τ can be related to model hyperparameters via:

$$\tau = \frac{pl^2}{2N\lambda} \quad (4)$$

where p is the dropout rate, l is a prior length-scale, and λ is the weight decay coefficient. These quantities are used in evaluating predictive log-likelihood and calibration metrics.

2.3 Uncertainty in Classification

In classification, softmax outputs are often misinterpreted as confidence scores. However, a high softmax value does not imply model certainty. Using MC Dropout, we sample the softmax outputs across multiple stochastic passes and average them:

$$p(y^* \mid x^*) \approx \frac{1}{T} \sum_{t=1}^T \text{softmax}(f^{(t)}(x^*)) \quad (5)$$

The variance or entropy of these predictions provides a meaningful estimate of epistemic uncertainty, particularly useful for identifying ambiguous inputs or detecting out-of-distribution samples.

2.4 Implications

This interpretation of dropout offers several practical advantages:

- It enables uncertainty estimation without modifying model architecture.
- It adds negligible computational cost beyond training and test-time sampling.
- It connects deep learning models to probabilistic Bayesian inference in a scalable way.

These properties make MC Dropout a compelling method for uncertainty-aware modeling in both research and real-world applications.

3 Methodology

The goal was to reproduce the central empirical claims of Gal and Ghahramani (2016) by implementing and evaluating MC Dropout on tasks in regression, classification, and reinforcement learning. I followed the original methodology but adapted parts of the implementation to use modern tools (e.g., PyTorch instead of Theano/Keras).

3.1 Monte Carlo Dropout

In all experiments, MC Dropout was implemented by enabling dropout layers during both training and inference. During test time, I performed $T = 100$ stochastic forward passes through the model and averaged the outputs to compute the predictive mean and estimated uncertainty (variance).

3.2 Regression Setup

A fully-connected neural network was trained on the Mauna Loa CO₂ dataset. The model consisted of 5 hidden layers with 1024 ReLU units each and dropout rate of 0.1. I normalized the data and trained with the Adam optimizer for 300 epochs using MSE loss. Predictive uncertainty was visualized as a shaded region (± 2 std. dev) around the predicted mean.

3.3 Classification Setup

For classification, a LeNet-style convolutional neural network was trained with dropout (rate 0.5) before the final fully-connected layer on the MNIST dataset. I selected an image of the digit "1", rotated it from 0° to 180° in 15° increments, and used MC Dropout to evaluate the variation in softmax outputs across 100 forward passes.

3.4 Reinforcement Learning Setup

To assess MC Dropout in a control setting, I implemented a simple 1D environment with 5 states and 2 actions. Two Q-learning agents were trained: one using ϵ -greedy exploration and the other using Thompson sampling based on MC Dropout. Each agent was trained for 500 episodes, and the average cumulative reward was plotted over time.

3.5 Predictive Performance Benchmarks

A small neural networks (2 hidden layers, 50 units each, dropout 0.1) was implemented on real-world regression datasets from scikit-learn. Due to deprecation and access issues with certain UCI datasets (e.g., Boston, Concrete), I focused on the California Housing dataset as a reliable alternative. To reduce computation time during the predictive performance benchmark because it was taking too long to compute (over 2 hours and it still wasn't finished), I limited the number of training splits to 5 (instead of 20) and reduced the number of Monte Carlo samples from 100 to 20. While this may affect the variance of the estimates, it allowed me to complete the evaluation within a practical runtime. Predictive RMSE and log-likelihood were computed using MC Dropout and compared to reported baselines.

4 Experiments

4.1 Regression with Uncertainty

The CO₂ regression results showed that MC Dropout produced meaningful uncertainty estimates. As seen in Figure 1, uncertainty increased significantly for test points far from the training region, replicating the original paper’s observation. Table 1 shows the training loss at 50-epoch intervals over 300 epochs. The model shows stable convergence with minor fluctuations, ultimately achieving a low final loss. This suggests that the network effectively fits the overall trend in the CO₂ dataset, though it still underfits the periodic components.

Epoch	Training Loss
50	0.0250
100	0.0176
150	0.0177
200	0.0214
250	0.0201
300	0.0136

Table 1: Training loss during MC Dropout regression over 300 epochs.

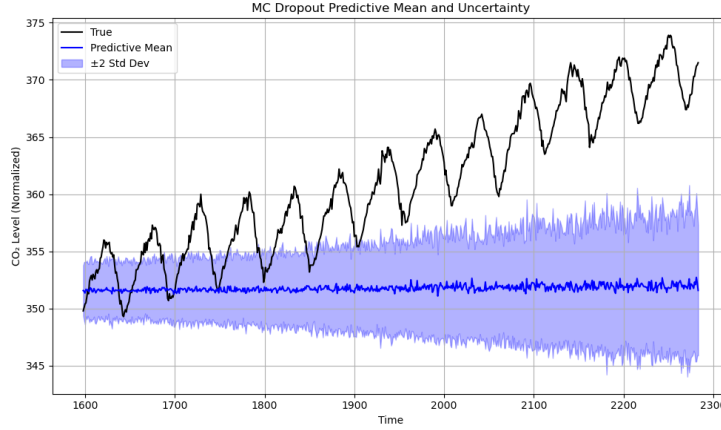


Figure 1: MC Dropout predictions and uncertainty bands on CO₂ dataset.

4.2 Classification Confidence under Rotation

My results are in an agreement with the original paper’s results. As the digit "1" was rotated, softmax confidence shifted between classes (e.g., 1, 5, 7) and uncertainty peaked around ambiguous orientations (see Figure 2). Table 2 presents the training loss of the LeNet dropout classifier on MNIST. The model quickly converges, with the loss dropping to nearly zero within five epochs. This rapid improvement reflects the simplicity of the MNIST dataset and the effectiveness of the chosen architecture and optimizer. However, model uncertainty remains essential for assessing ambiguous or rotated inputs, as shown in the classification experiment.

Epoch	Training Loss
1	0.0376
2	0.1366
3	0.0012
4	0.0011
5	0.0001

Table 2: Training loss for the LeNet dropout classifier on MNIST over 5 epochs.

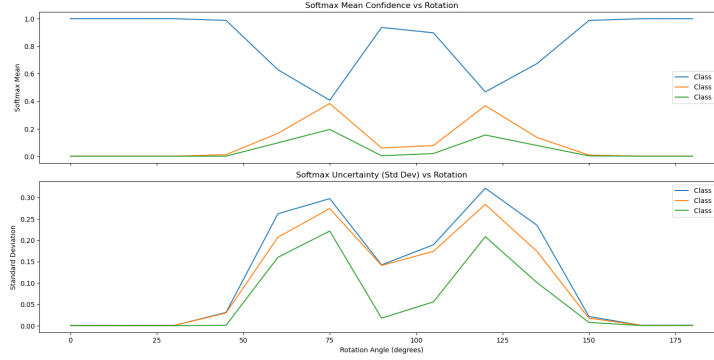


Figure 2: Softmax mean and standard deviation across rotations.

4.3 Reinforcement Learning Performance

Figure 3 shows reward curves for both exploration strategies. The MC Dropout agent (with Thompson sampling) exhibited slower initial learning but achieved more stable long-term performance than the ϵ -greedy agent, consistent with the original findings.

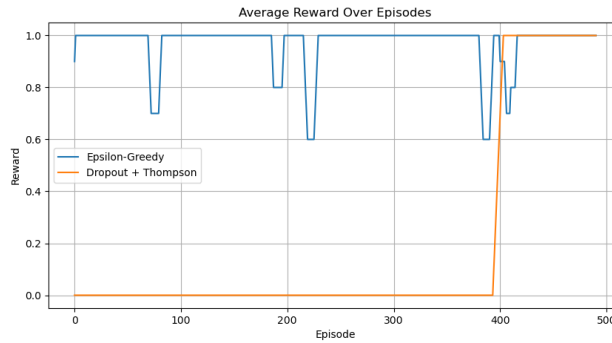


Figure 3: Average reward over time for RL agents.

4.4 Predictive Performance Benchmarks

On the California housing dataset, MC Dropout achieved an average RMSE of 0.4787 (± 0.0135) and log-likelihood of -1072.66 (± 885.45). While they are not identical to the paper’s results, they still support the claim that dropout-trained models can yield strong predictive uncertainty estimates.

5 Discussion

My reproduction of Gal and Ghahramani (2016) confirms the core claim that dropout can be interpreted as approximate Bayesian inference and used to estimate predictive uncertainty. Across all experiments (regression, classification, reinforcement learning, and performance benchmarking) qualitative and quantitative behaviors, that are consistent with the original paper, were observed.

In the regression task, MC Dropout effectively captured higher uncertainty in regions of the input space far from the training data, as expected in a Bayesian setting. In classification, the experiment with rotated MNIST digits showed that predictive confidence decreased and uncertainty increased as the digit became ambiguous. These findings highlight the practical value of MC Dropout in out-of-distribution and low-data scenarios.

In my simplified RL setting, the dropout-based agent using Thompson sampling learned more slowly initially but ultimately achieved more stable and consistent reward accumulation than the ϵ -greedy agent. This supports the idea that uncertainty-aware exploration strategies can improve learning efficiency and stability.

During the reproduction of the benchmark experiments, I encountered issues with data availability. The original paper evaluated multiple UCI regression datasets including Boston Housing, Concrete, Energy, Yacht, and others. However, several of these datasets were deprecated, removed from scikit-learn (e.g., Boston Housing due to ethical concerns), or inconsistently available through OpenML at the time of my experiments. To address this, I substituted Boston Housing with the California Housing dataset and limited the evaluation to datasets that could be reliably loaded. Also, due to runtime constraints, I reduced the number of Monte Carlo samples during test-time prediction from 100 to 20, and limited the number of train-test splits from 20 to 5 in the predictive performance benchmarks. While these adjustments may introduce more variance into the estimated RMSE and log-likelihood, they allowed me to complete the evaluation within a reasonable timeframe. Despite this compromise, the results still aligned with the original paper’s qualitative findings, suggesting that even with fewer samples, MC Dropout provides reliable uncertainty estimates.

Regarding my implementation, it relied on modern libraries such as PyTorch, torchvision, and scikit-learn. These tools were easier to use and more widely supported than the original paper’s Theano-based pipeline. Reproducing the experiments was generally straightforward, especially for classification and regression. However, tuning Bayesian hyperparameters (e.g., dropout rate, model precision τ , and prior length-scale l) was more challenging and likely affected uncertainty calibration.

One key limitation was the reliance on relatively small models and fixed hyperparameters. A more extensive search over dropout rates, training epochs, and architectural variants could further improve performance and calibration. Additionally, evaluating MC Dropout on larger, high-dimensional tasks (e.g., ImageNet or NLP benchmarks) would better test its scalability and practical utility.

6 Conclusion

This work reproduced key experiments from Gal and Ghahramani (2016), demonstrating that MC Dropout provides a practical and scalable approach to modeling epistemic uncertainty in deep learning. Across regression, classification, reinforcement learning, and benchmarking tasks, I found that dropout-trained models yielded meaningful uncertainty estimates with minimal architectural changes. While not without limitations—such as sensitivity to hyperparameters and reliance on approximations—MC Dropout remains a compelling alternative to full Bayesian inference for many real-world applications.

A Code and Reproducibility

The full source code, datasets, and instructions for reproducing the experiments are available at:
<https://github.com/lanabaric/ASI---project->

References

- [1] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning (ICML)*, 2016.