

Dropout as a Bayesian Approximation

appliedAI Seminar — Uncertainty Quantification in Deep Learning

Maternus Herold

20.10.2022



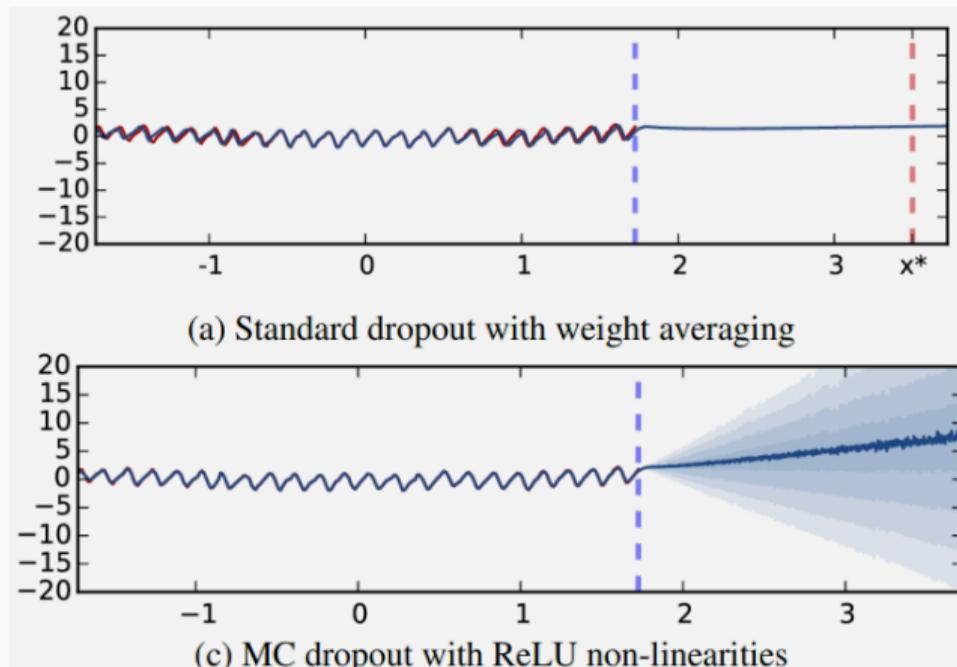
**INITIATIVE FOR
APPLIED ARTIFICIAL
INTELLIGENCE**

Agenda

1. Motivation
2. Quick intro to Uncertainty Quantification
3. Recap of Dropout in Neural Networks
4. MC Dropout for Uncertainty Quantification in Deep Learning

Motivation

Why do we want to quantify the uncertainty in our models?



Figures showing the true value (red) and the predictive mean (blue) including $\pm 2\sigma$ for a model evaluated on the Mauna Loa CO₂ concentration dataset [4]. The blue shades indicate half a std. deviation each.

Uncertainty Quantification

Uncertainty is divided into aleatoric and epistemic uncertainty

Uncertainty is divided into aleatoric and epistemic uncertainty

Aleatoric uncertainty

- Stochastic uncertainty or unknowns, responsible for different outcomes of the same experiment
- Often due to insufficient measurement quality
- Does not preclude the existence of the “unknowns”
- Uncertainty trapped “in the data”
- Methods: Quantile regression, mixture of densities, max. likelihood est.

Uncertainty is divided into aleatoric and epistemic uncertainty

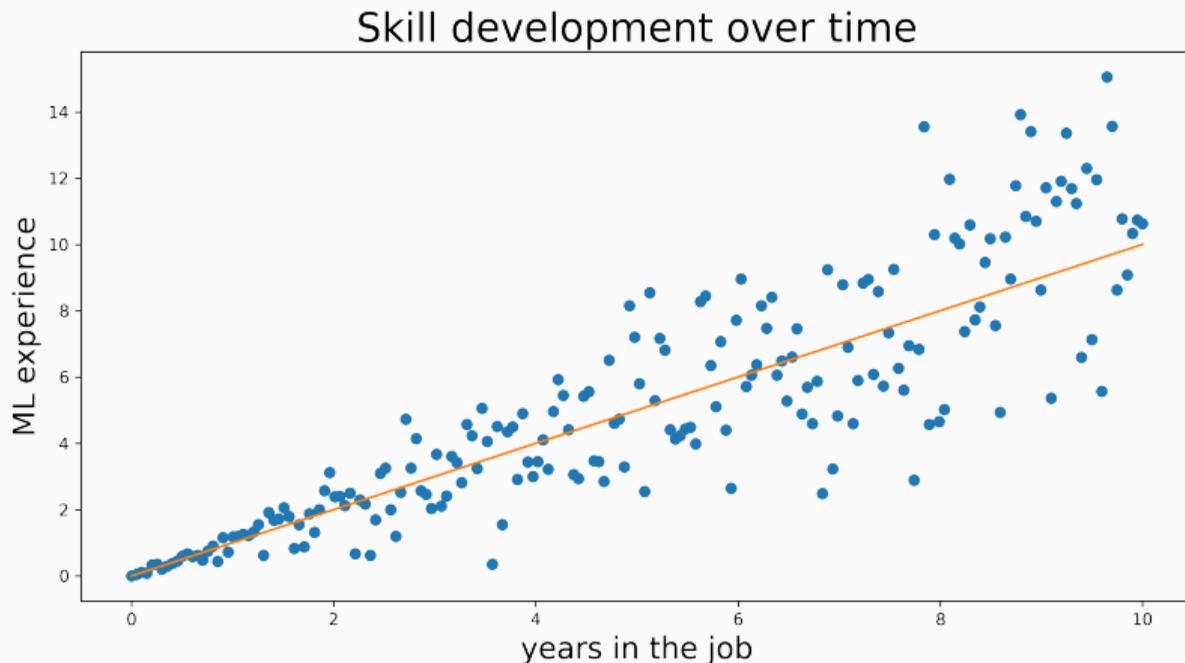
Aleatoric uncertainty

- Stochastic uncertainty or unknowns, responsible for different outcomes of the same experiment
- Often due to insufficient measurement quality
- Does not preclude the existence of the “unknowns”
- Uncertainty trapped “in the data”
- Methods: Quantile regression, mixture of densities, max. likelihood est.

Epistemic uncertainty

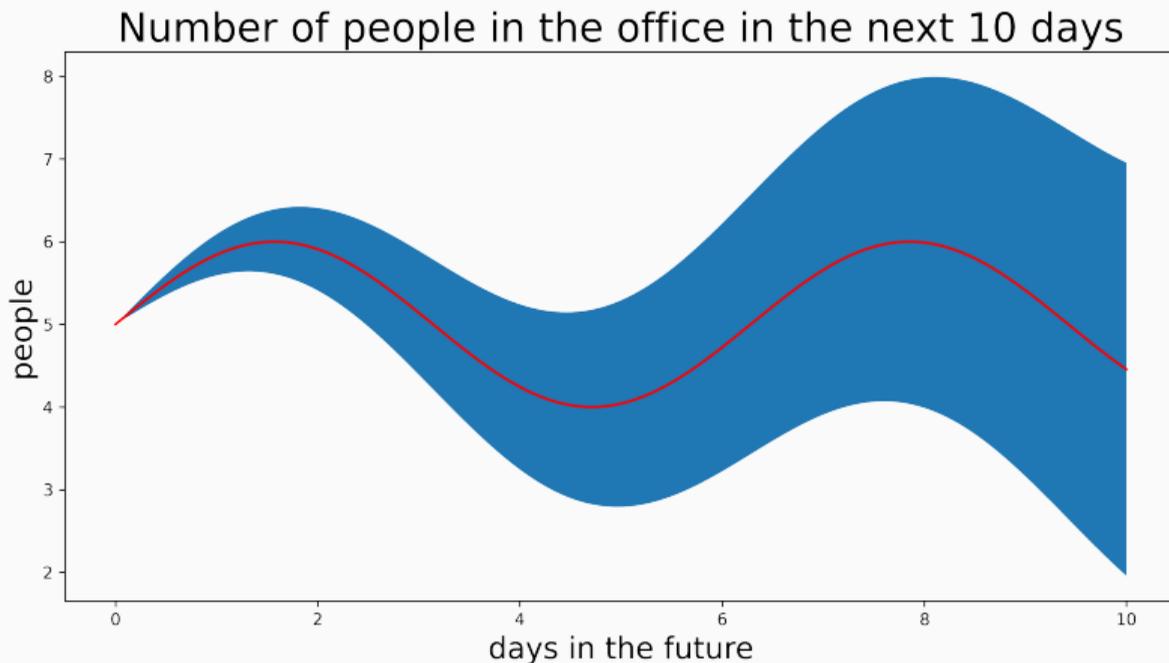
- Systematic uncertainty, due to possibly available information but practically not used
- Inaccuracies in measurements, unused data / features, neglecting effects
- Captured by the model due to design decisions, data used, etc.
- Methods: Bayesian Neural Networks, MC dropout or deep ensembles

Illustration of aleatoric uncertainty



Based on the current information (years in the job), the increasing variance in ML experience cannot be explained.

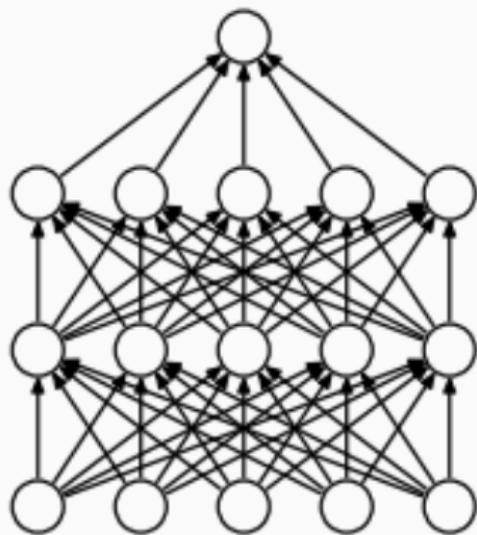
Illustration of epistemic uncertainty



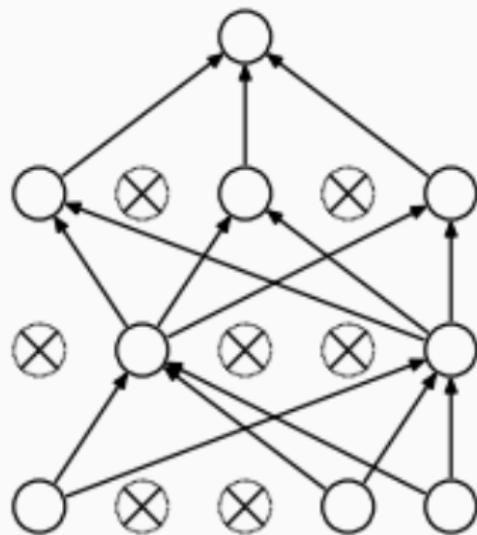
While the model returns point estimates, it is uncertain about its predictions, especially further in the future.

Recap of Dropout

Dropping our neurons to prevent overfitting



(a) Standard Neural Net



(b) After applying dropout.

Illustration of a standard neural network vs. one with dropping out nodes [6].

Model definition including dropping out nodes

$$\begin{aligned} z_i^{l+1} &= w_i^{l+1} \cdot y^l + b_i^{l+1} \\ y_i^{l+1} &= \sigma(z_i^{l+1}) \end{aligned} \quad (1)$$

$$\begin{aligned} \eta_j^l &\sim \mathcal{B}(p) \\ \tilde{y}^l &= \eta^l \cdot y^l \\ z_i^{l+1} &= w_i^{l+1} \cdot \tilde{y}^l + b_i^{l+1} \\ y_i^{l+1} &= \sigma(z_i^{l+1}) \end{aligned} \quad (2)$$

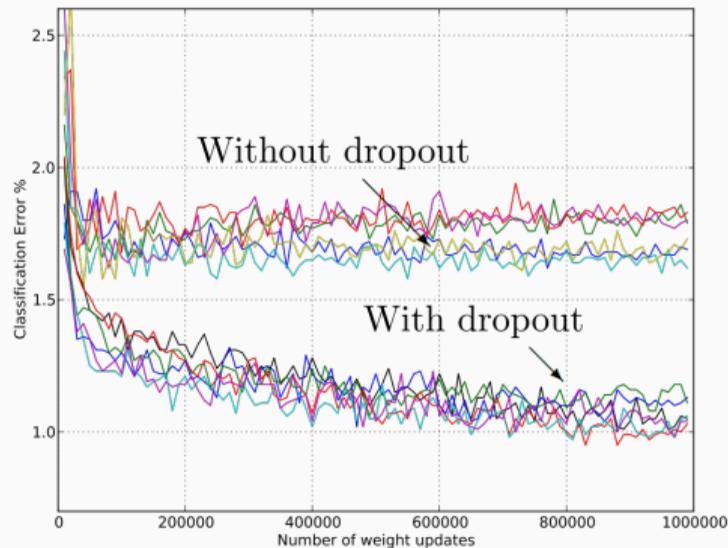
Where the keep probability is p and layers are denoted by $l \in \{0, \dots, L - 1\}$ [6].

Dropout is primarily used to prevent overfitting

- Reduce risk of overfitting
- Model ensemble without optimizing hundreds of models
- Activate all parts of the model

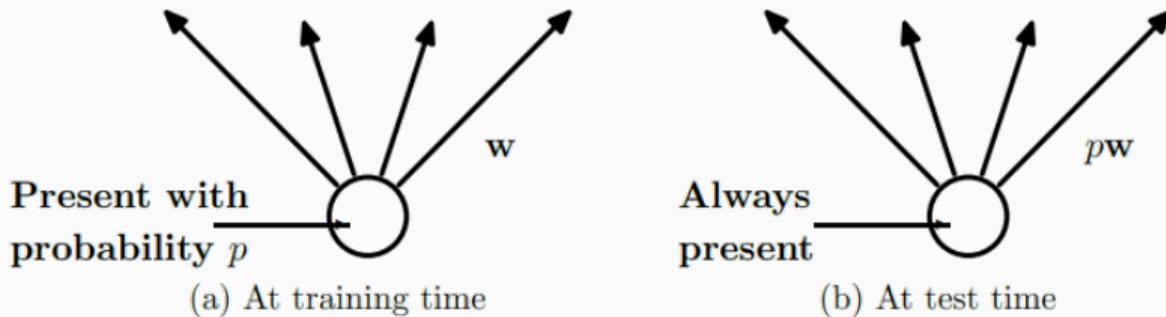
Dropout is primarily used to prevent overfitting

- Reduce risk of overfitting
- Model ensemble without optimizing hundreds of models
- Activate all parts of the model



Dropout improves performance on the test data for different architectures and configurations [6].

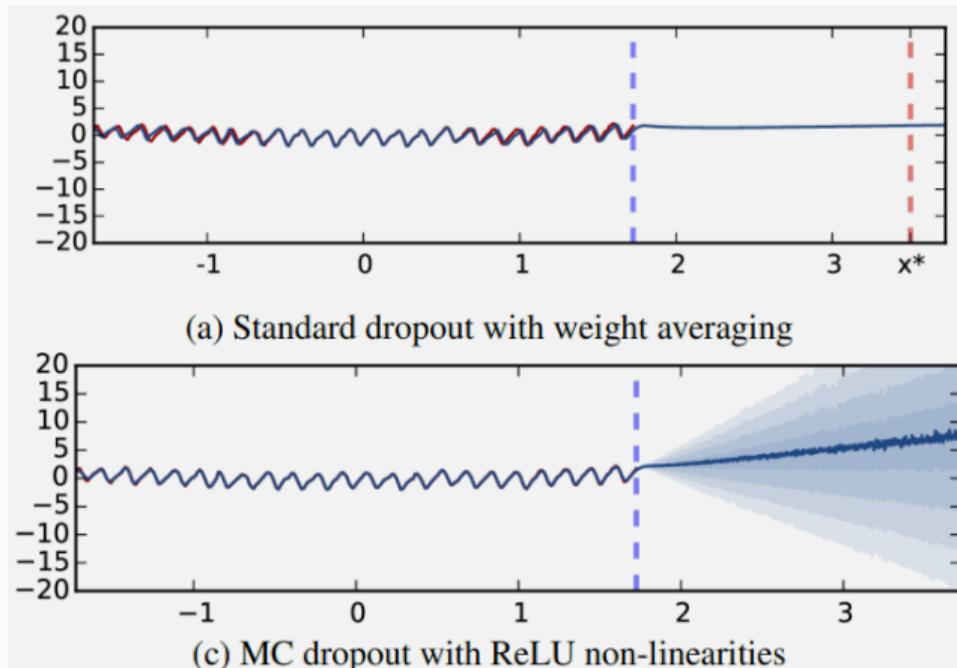
Dropout is switched off during inference — all nodes are active



As dropping nodes is turned off during inference, the weights need to be scaled to not get larger than expected [6].

Monte Carlo Dropout

Q: How do we get the model's uncertainty?



Figures showing the true value (red) and the predictive mean (blue) including $\pm 2\sigma$ for a model evaluated on the Mauna Loa CO₂ concentration dataset [4]. The blue shades indicate half a std. deviation each.

Retrieving the predictive mean and variance

Given a model \mathcal{M} which was trained with dropout and a *keep probability* of p .
Generate $\{(\hat{y}_i, \mathbf{x})\}_{i=1}^N$ by **keeping dropout on**.

Estimators for $\mathbb{E}(\mathbf{y})$, $\text{Var}(\mathbf{y})$ read

$$\begin{aligned}\mathbb{E}(\mathbf{y}) &\approx \frac{1}{N} \sum_{i=1}^N \hat{y}_i(\mathbf{x}) \\ \text{Var}(\mathbf{y}) &\approx \tau^{-1} \mathbf{I} + \frac{1}{N} \sum_{i=1}^N \hat{y}_i(\mathbf{x}) \hat{y}_i(\mathbf{x})^T - \mathbb{E}(\mathbf{y}) \mathbb{E}(\mathbf{y})^T\end{aligned}\tag{3}$$

Where $\tau = \frac{l^2 p}{2N\lambda}$ the GP precision; τ, p obtained via hyperparameter search [5].

Algorithm 1: Implementation of Monte Carlo Dropout by using dropout during inference as well and computing statistics of the samples from the predictive distribution.

```
1:  $\mathcal{M} \leftarrow \mathcal{M}.\text{train}()$   
2:  $\hat{y}_i \leftarrow \mathcal{M}(x) \forall i \in [1, \dots, N]$   
3:  
4:  $\hat{\mathbb{E}}(y) \leftarrow \hat{y}.\text{mean}()$   
5:  $\hat{\sigma}(y) \leftarrow \hat{y}.\text{std}()$   
6:  
7: return  $\hat{\mathbb{E}}(y), \hat{\sigma}(y)$ 
```

Why does it work?

💡 *“Averaging over forward passes is equivalent to Monte Carlo integration over a Gaussian Process posterior approximation.”*

→ We'll familiarize with the idea **but** I refer to [3, 4] for a detailed discussion.

Why does it work?

💡 “Averaging over forward passes is equivalent to Monte Carlo integration over a Gaussian Process posterior approximation.”

→ We'll familiarize with the idea **but** I refer to [3, 4] for a detailed discussion.

Target: predictive distribution for a new observation \mathbf{x}^*

$$\begin{aligned} p(y^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) &= \int p(y^*|\mathbf{x}^*, \omega) p(\omega|\mathbf{X}, \mathbf{Y}) d\omega \\ &\approx \int p(y^*|\mathbf{x}^*, \omega) q_{\theta^*}(\omega) d\omega, \quad \theta^* = \underset{\theta}{\operatorname{argmin}} \operatorname{KL}(q_{\theta}(\omega)|p(\omega|\mathbf{X}, \mathbf{Y})) \end{aligned} \tag{4}$$

Finding $q_{\theta^*} \Rightarrow$ Minimizing the KL-Divergence

💡 Minimizing KL-Div. \Leftrightarrow maximizing log evidence lower bound [1].

$$\begin{aligned}\log p(\mathbf{Y}|\mathbf{X}) &= \log \int p(\mathbf{Y}|\mathbf{X}, \omega) p(\omega) d\omega \\ &= \log \int p(\mathbf{Y}|\mathbf{X}, \omega) p(\omega) \frac{q_{\theta}(\omega)}{q_{\theta}(\omega)} d\omega \\ &= \log \left(\mathbb{E}_{q_{\theta}} \frac{p(\mathbf{Y}|\mathbf{X}, \omega) p(\omega)}{q_{\theta}(\omega)} \right) \\ &\stackrel{\text{Jensen's ineq.}}{\geq} \mathbb{E}_{q_{\theta}} \log \left(\frac{p(\mathbf{Y}|\mathbf{X}, \omega) p(\omega)}{q_{\theta}(\omega)} \right) \\ L_{\text{VI}} &= \int q_{\theta}(\omega) \log(p(\mathbf{Y}|\mathbf{X}, \omega)) d\omega - \text{KL}(q_{\theta}|p(\omega))\end{aligned}\tag{5}$$

Maximizing the log evidence lower bound

Using MC integration to approximate the integral with one sample $\omega \sim q_\theta$ and optimizing w.r.t. θ :

$$\hat{L}(\theta) = \log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) - \text{KL}(q_\theta|p(\omega)) \quad (6)$$

Connecting the VI approach with dropout

$\Rightarrow q_{\theta}(\omega)$ is a distribution over weights of a network. Thus, define $q_{\theta}(\omega)$ s.t. weights can be turned off.

$$q_{\theta^*} = \prod_i q_{M_i}(W_i), \quad i \in [1, \dots, L]$$

$$q_{M_i}(W_i) = M_i \text{diag}(z_{i,j}) \tag{7}$$

$$z_{i,j} \sim \mathcal{B}(p_i)$$

$$W_i \sim q_{M_i}(W_i)$$

- Requires large number of samples making inference slow ($> 10^3$ in [4])
- Inference time is crucial in a lot of applications

Single-shot Monte Carlo Dropout

- “Succession” of the original Monte Carlo Dropout
- Preserves advantages of Bayesian NN without being slower than non-Bayesian NN [2]
- Idea: approximate the expected value and variance of the MC Dropout per layer

Approximation via Moment Propagation

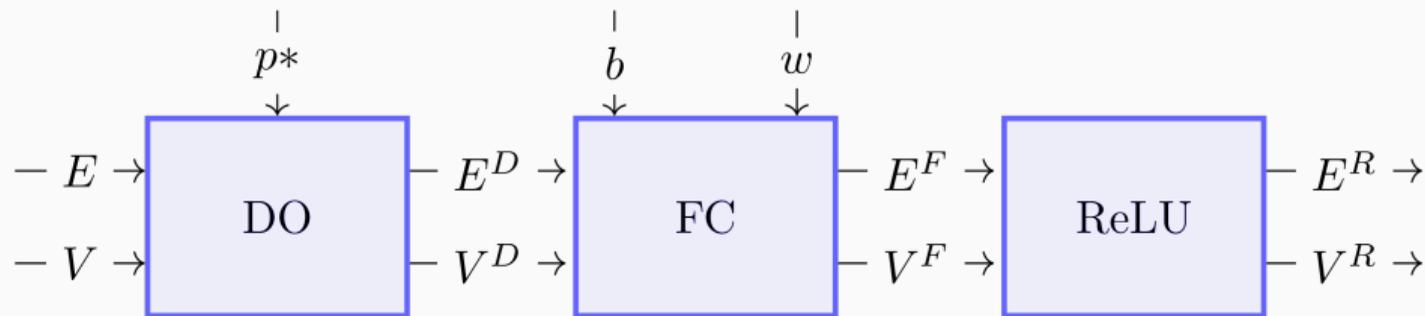


Illustration of Moment Propagation through dropout, fully-connected and ReLU layers [2]. The expectation and variance *flow* through the network during a single forward pass

Dropout Layer per layer with nodes i

$$E^i = E_{\text{in}}^i \cdot p$$

$$V^i \stackrel{X \parallel Y}{=} V_{\text{in}}^i \cdot p(1-p) + V_{\text{in}}^i p^2 + (E_{\text{in}}^i)^2 \cdot p(1-p) \quad (8)$$

Approximation via Moment Propagation ctd.

Dropout Layer per layer with nodes i

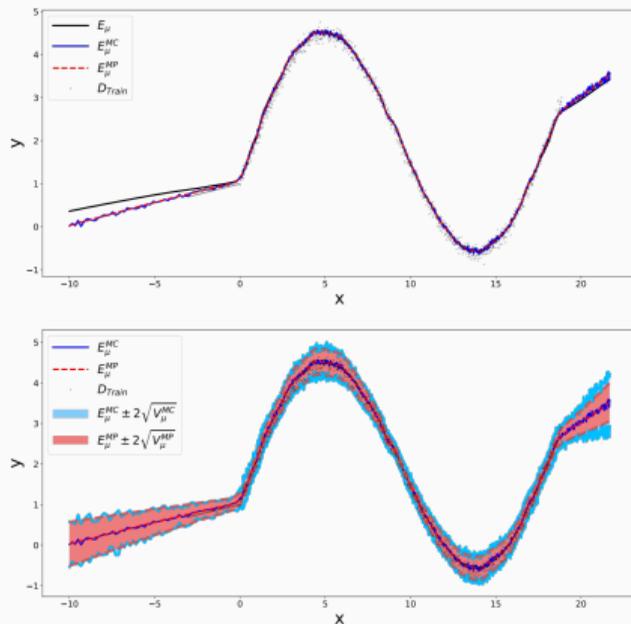
$$E^i = E_{\text{in}}^i \cdot p$$

$$V^i \stackrel{X \parallel Y}{=} V_{\text{in}}^i \cdot p(1-p) + V_{\text{in}}^i p^2 + (E_{\text{in}}^i)^2 \cdot p(1-p) \quad (8)$$

Dense Layer

$$E^i = \sum_j^{N_{\text{in}}} w_{ji} E_{\text{in}}^j + b^i$$
$$V^i = \sum_j^{N_{\text{in}}} w_{ji}^2 V_{\text{in}}^j \quad (9)$$

Comparison of Single Shot MC and MC Dropout



Visual comparison of the two methods using MC Dropout. The models were trained on data in the range of $[-3, 19]$. In the upper figure, the predictions are shown. In the lower panel, the uncertainty is included, showing similar results [2].

Comparison of Single Shot MC and MC Dropout ctd.

DATASET	N	Q	TEST RMSE		TEST NLL		TEST RT [s]	
			MC	MP	MC	MP	MC	MP
BOSTON	506	13	3.14 \pm 0.20	3.10 \pm 0.20	2.57 \pm 0.07	2.56 \pm 0.08	2.51 \pm 0.03	0.04 \pm 0.00
CONCRETE	1,030	8	5.46 \pm 0.12	5.40 \pm 0.12	3.12 \pm 0.02	3.13 \pm 0.03	3.37 \pm 0.04	0.04 \pm 0.00
ENERGY	768	8	1.65 \pm 0.05	1.61 \pm 0.05	1.95 \pm 0.04	2.01 \pm 0.04	2.84 \pm 0.03	0.04 \pm 0.00
KIN8NM	8,192	8	0.08 \pm 0.00	0.08 \pm 0.00	-1.10 \pm 0.01	-1.11 \pm 0.01	7.37 \pm 0.06	0.04 \pm 0.00
NAVAL	11,934	16	0.00 \pm 0.00	0.00 \pm 0.00	-4.36 \pm 0.01	-3.60 \pm 0.01	9.69 \pm 0.11	0.04 \pm 0.00
POWER	9,568	4	4.05 \pm 0.04	4.04 \pm 0.04	2.82 \pm 0.01	2.84 \pm 0.01	6.85 \pm 0.07	0.04 \pm 0.00
PROTEIN	45,730	9	4.42 \pm 0.03	4.41 \pm 0.02	2.90 \pm 0.00	2.91 \pm 0.00	31.38 \pm 0.09	0.05 \pm 0.00
WINE	1,599	11	0.63 \pm 0.01	0.63 \pm 0.01	0.95 \pm 0.01	0.95 \pm 0.01	4.78 \pm 0.01	0.04 \pm 0.00
YACHT	308	6	2.93 \pm 0.22	2.91 \pm 0.26	2.35 \pm 0.07	2.11 \pm 0.07	2.01 \pm 0.01	0.04 \pm 0.00

Benchmarking both approaches on the UCI benchmarking datasets where $N = 10^4$ forward passes were used for the MC Dropout model. Both methods show similar results on the test data while the approach of Moment Propagation has a significant advantage w.r.t. inference time [2].

Dropout as a Bayesian Approximation

appliedAI Seminar — Uncertainty Quantification in Deep Learning

Maternus Herold

20.10.2022



**INITIATIVE FOR
APPLIED ARTIFICIAL
INTELLIGENCE**

Literatur

- [1] Christopher M Bishop und Nasser M Nasrabadi. *Pattern recognition and machine learning*. Bd. 4. 4. Springer, 2006.
- [2] Kai Brach, Beate Sick und Oliver Dürr. *Single Shot MC Dropout Approximation*. 2020. URL: <https://arxiv.org/abs/2007.03293>.
- [3] Yarin Gal und Zoubin Ghahramani. *Dropout as a Bayesian Approximation: Appendix*. 2015. URL: <https://arxiv.org/abs/1506.02157>.

- [4] Yarin Gal und Zoubin Ghahramani. “Dropout as a bayesian approximation: Representing model uncertainty in deep learning”. In: *international conference on machine learning*. PMLR. 2016, S. 1050–1059.
- [5] Yarin Gal und Jishnu Mukhoti. *Code Examples for Dropout as Bayesian Approximation*.
<https://github.com/yaringal/DropoutUncertaintyExps>. [Online; accessed 19-Oct-2022]. 2016.
- [6] Nitish Srivastava u. a. “Dropout: A Simple Way to Prevent Neural Networks from Overfitting”. In: *Journal of Machine Learning Research* 15.56 (2014), S. 1929–1958. URL:
<http://jmlr.org/papers/v15/srivastava14a.html>.