

Well-Calibrated Regression Uncertainty in Medical Imaging with Deep Learning

Max-Heinrich Laves, Sontje Ihler, Jacob F. Fast, Lüder A. Kahrs,
Tobias Ortmaier

Medical Imaging With Deep Learning (MIDL)

6–9 July 2020

Regression in Medical Imaging

- 1 Age estimation from hand CT (Halabi et al., 2019)
- 2 Natural landmark localization (Payer et al., 2019)
- 3 Cell detection in histology (Xie et al., 2018)
- 4 Instrument pose estimation (Gessert et al., 2018)
- 5 Deformable registration (Dalca et al., 2019)

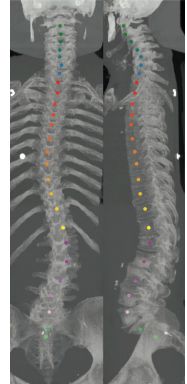
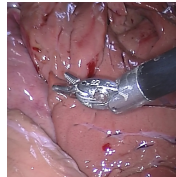


Figure: Medical regression tasks.

Predictive Uncertainty

- Reliable predictions are **crucial**
- Two types of uncertainty (Kendall et al., 2017)



Figure: Different types of uncertainty, note object boundaries (Kendall et al., 2017).

Predictive Uncertainty

- Reliable predictions are **crucial**
- Two types of uncertainty (Kendall et al., 2017)
- Aleatoric
 - Arises from data directly (e. g. sensor noise)

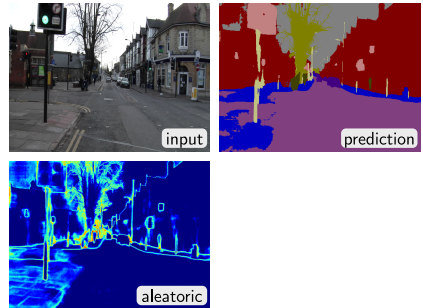


Figure: Different types of uncertainty, note object boundaries (Kendall et al., 2017).

Predictive Uncertainty

- Reliable predictions are **crucial**
- Two types of uncertainty (Kendall et al., 2017)
- Aleatoric
 - Arises from data directly (e. g. sensor noise)
- Epistemic
 - From limited training data

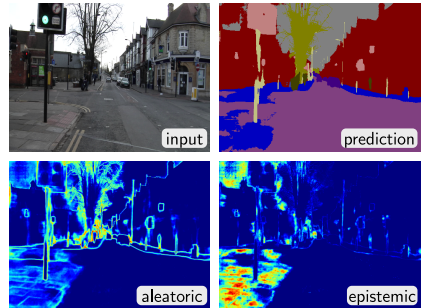


Figure: Different types of uncertainty, note object boundaries (Kendall et al., 2017).

Predictive Uncertainty

- Reliable predictions are **crucial**
- Two types of uncertainty (Kendall et al., 2017)
- Aleatoric
 - Arises from data directly (e. g. sensor noise)
- Epistemic
 - From limited training data
- **Bayesian Neural Networks**

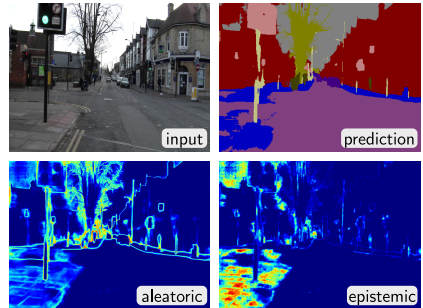


Figure: Different types of uncertainty, note object boundaries (Kendall et al., 2017).

Predictive Uncertainty

- Reliable predictions are **crucial**
- Two types of uncertainty (Kendall et al., 2017)
- Aleatoric
 - Arises from data directly (e. g. sensor noise)
- Epistemic
 - From limited training data
- **Bayesian Neural Networks**
- Uncertainty is **miscalibrated**

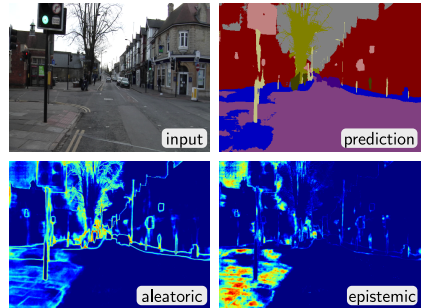


Figure: Different types of uncertainty, note object boundaries (Kendall et al., 2017).

Estimation of Aleatoric Uncertainty I

Conditional Log-Likelihood for Regression

$$\mathbf{f}_{\theta}(\mathbf{x}) = [\hat{\mathbf{y}}(\mathbf{x}), \hat{\sigma}^2(\mathbf{x})], \quad \hat{\mathbf{y}} \in \mathbb{R}^d$$

$$\mathcal{L}(\theta) = \sum_{i=1}^m \frac{1}{\hat{\sigma}^2(\mathbf{x}_i)} \|\mathbf{y}_i - \hat{\mathbf{y}}(\mathbf{x}_i)\|^2 + \log \hat{\sigma}^2(\mathbf{x}_i)$$

Estimation of Aleatoric Uncertainty I

Conditional Log-Likelihood for Regression

$$\mathbf{f}_{\theta}(\mathbf{x}) = [\hat{\mathbf{y}}(\mathbf{x}), \hat{\sigma}^2(\mathbf{x})], \quad \hat{\mathbf{y}} \in \mathbb{R}^d$$

$$\mathcal{L}(\theta) = \sum_{i=1}^m \frac{1}{\hat{\sigma}^2(\mathbf{x}_i)} \|\mathbf{y}_i - \hat{\mathbf{y}}(\mathbf{x}_i)\|^2 + \log \hat{\sigma}^2(\mathbf{x}_i)$$

Problem Statement

Minimizing NLL w.r.t. $\hat{\sigma}^2(\mathbf{x}_i)$ yields

$$\hat{\sigma}^2(\mathbf{x}_i) = \arg \min_{\hat{\sigma}^2(\mathbf{x}_i)} \mathcal{L} = \|\mathbf{y}_i - \hat{\mathbf{y}}(\mathbf{x}_i)\|^2 \quad \forall i.$$

Estimation of Aleatoric Uncertainty II

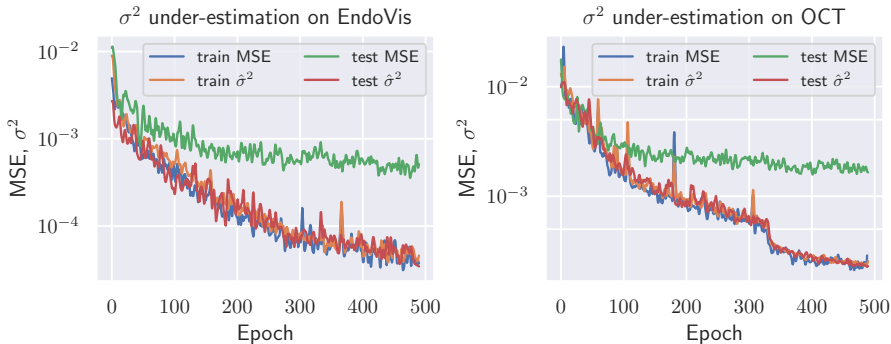


Figure: σ^2 is estimated relative to the MSE.

Recalibration of Standard Deviation

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}(\mathbf{x}), (s \cdot \hat{\sigma})^2(\mathbf{x}))$$

Recalibration of Standard Deviation

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}(\mathbf{x}), (s \cdot \hat{\sigma})^2(\mathbf{x}))$$

$$\mathcal{L}(s) = m \log(s) + \frac{s^{-2}}{2} \sum_{i=1}^m \hat{\sigma}^2(\mathbf{x}_i) \|\mathbf{y}^{(i)} - \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}}^{(i)}\|^2$$

Recalibration of Standard Deviation

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}(\mathbf{x}), (s \cdot \hat{\sigma})^2(\mathbf{x}))$$

$$\mathcal{L}(s) = m \log(s) + \frac{s^{-2}}{2} \sum_{i=1}^m \hat{\sigma}^2(\mathbf{x}_i) \|\mathbf{y}^{(i)} - \hat{\mu}_{\theta}^{(i)}\|^2$$

$$s = \pm \sqrt{\frac{1}{m} \sum_{i=1}^m (\hat{\sigma}_{\theta}^{(i)})^{-2} \|\mathbf{y}^{(i)} - \hat{\mathbf{y}}_{\theta}^{(i)}\|^2}.$$

→ We refer to this as σ scaling.

Well-Calibrated Estimation of Predictive Uncertainty

- So far: maximum posterior point estimate $\hat{\theta}$
- Bayesian model with Monte Carlo dropout VI (Gal et al., 2016)

Well-Calibrated Estimation of Predictive Uncertainty

- So far: maximum posterior point estimate $\hat{\theta}$
- Bayesian model with Monte Carlo dropout VI (Gal et al., 2016)

Predictive Uncertainty

Combines aleatoric (data) and epistemic (model) uncertainty (Kendall et al., 2017).

$$\hat{\Sigma}^2 = \underbrace{\frac{1}{N} \sum_{n=1}^N \left(\hat{y}_n - \frac{1}{N} \sum_{n=1}^N \hat{y}_n \right)^2}_{\text{epistemic}} + \underbrace{\frac{1}{N} \sum_{n=1}^N \hat{\sigma}_n^2}_{\text{aleatoric}}$$

VI under-estimates predictive variance.

→ Apply σ scaling to calibrate predictive uncertainty $(s \cdot \hat{\Sigma}(\mathbf{x}))^2$.

Definition of Miscalibration

Difference in expectation between predictive error and uncertainty

$$\mathbb{E}_{\hat{\Sigma}^2} \left[|(\|\mathbf{y} - \hat{\mathbf{y}}\|^2 \mid \hat{\Sigma}^2 = \Sigma^2) - \Sigma^2| \right] \quad \forall \{ \Sigma^2 \in \mathbb{R} \mid \Sigma^2 \geq 0 \} \quad (1)$$

Quantification of Miscalibration

Definition of Miscalibration

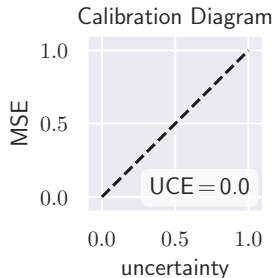
Difference in expectation between predictive error and uncertainty

$$\mathbb{E}_{\hat{\Sigma}^2} \left[\left| (\|\mathbf{y} - \hat{\mathbf{y}}\|^2 \mid \hat{\Sigma}^2 = \Sigma^2) - \Sigma^2 \right| \right] \quad \forall \{ \Sigma^2 \in \mathbb{R} \mid \Sigma^2 \geq 0 \} \quad (1)$$

Uncertainty Calibration Error

Partitioning into M bins (Guo et al., 2017)

$$\text{UCE} := \sum_{m=1}^M \frac{|B_m|}{n} |\text{err}(B_m) - \text{uncert}(B_m)|$$



- Four medical datasets with $\mathbf{y} \in \mathbb{R}^d$
 - ① tumor cellularity in breast histology ($d = 1$) (Martel et al., 2019)
 - ② RNSA bone age data set ($d = 1$) (Halabi et al., 2019)
 - ③ EndoVis surgical instrument tracking ($d = 2$) (EndoVis, 2015)
 - ④ needle pose estimation from 3D-OCT, own dataset ($d = 6$)¹
- Uncertainty calibration
- Rejection of uncertain predictions
- Out-of-distribution detection (see paper)

¹github.com/mlaves/3doct-pose-dataset

Intra-Training Calibration

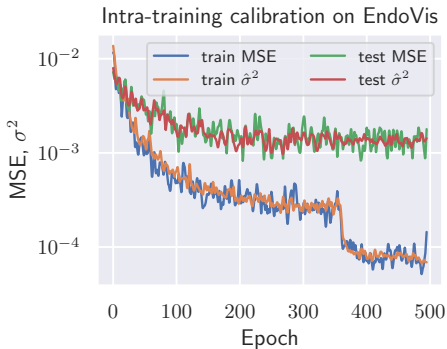


Figure: $\hat{\sigma}^2$ is **not** under-estimated.

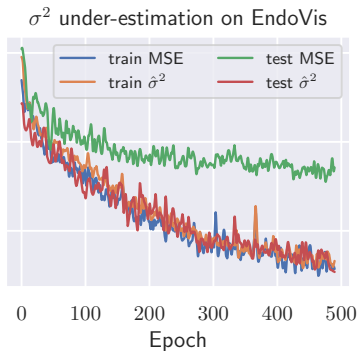
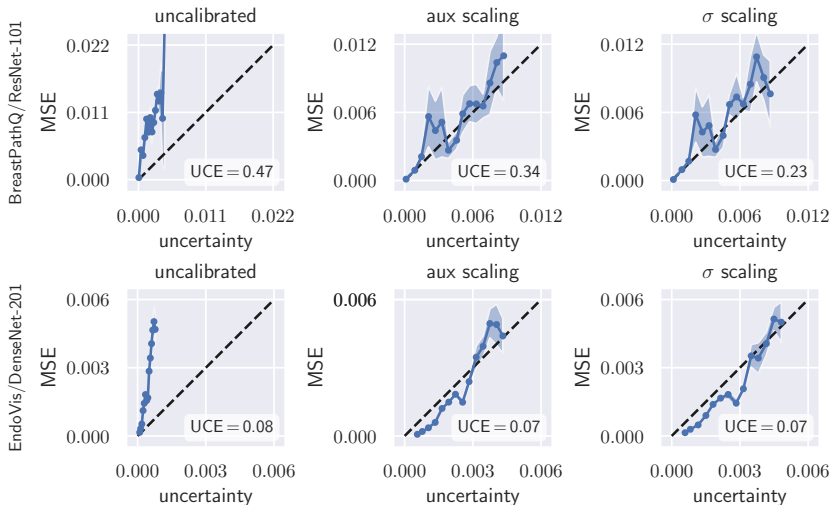


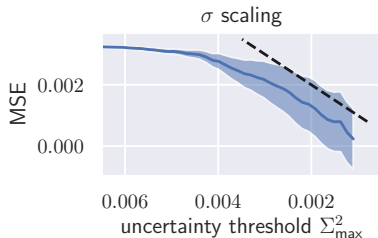
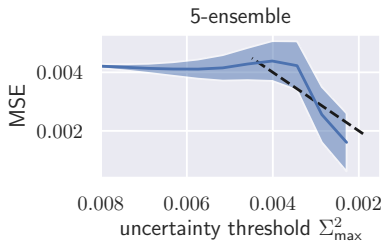
Figure: $\hat{\sigma}^2$ is under-estimated.

Calibration Diagrams



Rejection Experiments

- Uncertainty threshold Σ_{\max}^2
- Reject, where $\hat{\Sigma}^2 > \Sigma_{\max}^2$
- Reduce Σ_{\max}^2 , observe test MSE
- Compare to ensemble uncertainty
- σ scaling: **monotonic** decrease



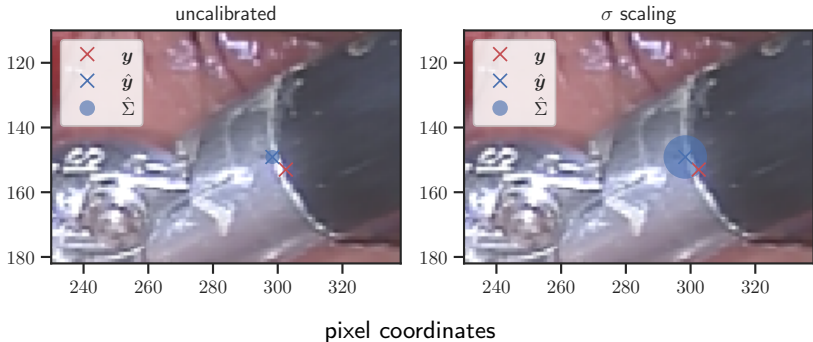


Figure: After σ scaling, the uncertainty better reflects the predictive error.

- Well-calibrated predictive uncertainty for regression
- Miscalibration is considerably reduced
- If already calibrated: $s \rightarrow 1$
- Reliably detects distribution shift
- Ensemble outperformed on rejection task

- Well-calibrated predictive uncertainty for regression
- Miscalibration is considerably reduced
- If already calibrated: $s \rightarrow 1$
- Reliably detects distribution shift
- Ensemble outperformed on rejection task
- Simple to implement
- Does not affect accuracy
- Closes gap between test MSE and uncertainty
- Well-calibrated uncertainty should be considered in any medical imaging task with deep learning

Well-Calibrated Regression Uncertainty in Medical Imaging with Deep Learning

Max-Heinrich Laves, Sontje Ihler, Jacob F. Fast, Lüder A. Kahrs,
Tobias Ortmaier

Medical Imaging With Deep Learning (MIDL)

6–9 July 2020

- Dalca, Adrian V. et al. (2019). “Unsupervised learning of probabilistic diffeomorphic registration for images and surfaces”. In: *Med Image Anal* 57, pp. 226–236. DOI: <https://doi.org/10.1016/j.media.2019.07.006>.
- EndoVis (2015). *Instrument Subchallenge Dataset*. <https://opencas.webarchiv.kit.edu/?q=node/30>.
- Gal, Yarin and Zoubin Ghahramani (2016). “Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning”. In: *ICML*, pp. 1050–1059.
- Gessert, Nils, Matthias Schlüter, and Alexander Schlaefher (2018). “A deep learning approach for pose estimation from volumetric OCT data”. In: *Med Image Anal* 46, pp. 162–179. DOI: [10.1016/j.media.2018.03.002](https://doi.org/10.1016/j.media.2018.03.002).
- Guo, Chuan et al. (2017). “On Calibration of Modern Neural Networks”. In: *ICML*, pp. 1321–1330.

- Halabi, Safwan S. et al. (2019). “The RSNA Pediatric Bone Age Machine Learning Challenge”. In: *Radiol* 290.2, pp. 498–503. DOI: [10.1148/radiol.2018180736](https://doi.org/10.1148/radiol.2018180736).
- Kendall, Alex and Yarin Gal (2017). “What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?” In: *NeurIPS*, pp. 5574–5584.
- Martel, A. L. et al. (2019). “Assessment of Residual Breast Cancer Cellularity after Neoadjuvant Chemotherapy using Digital Pathology [Data set]”. In: *The Cancer Imaging Archive*. DOI: [10.7937/TCIA.2019.4YIBTJNO](https://doi.org/10.7937/TCIA.2019.4YIBTJNO).
- Payer, Christian et al. (2019). “Integrating spatial configuration into heatmap regression based CNNs for landmark localization”. In: *Med Image Anal* 54, pp. 207–219. DOI: [10.1016/j.media.2019.03.007](https://doi.org/10.1016/j.media.2019.03.007).
- Xie, Yuanpu et al. (2018). “Efficient and robust cell detection: A structured regression approach”. In: *Med Image Anal* 44, pp. 245–254. DOI: [10.1016/j.media.2017.07.003](https://doi.org/10.1016/j.media.2017.07.003).