

# Uncertainty in Deep Learning

Kevin Barkevich





# Uncertainty: What is it?

## Aleatoric:

- “*Alea*” - Latin for “dice”
- Represents the variability/randomness in the outcome of an experiment
- High uncertainty indicates the presence of meaningless “random noise”
- **Not** reducible with more training data

## Epistemic:

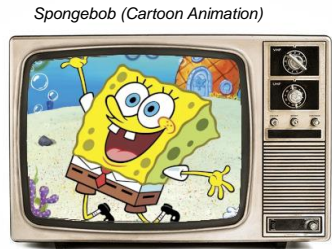
- Represents a lack of knowledge in what the experiment’s outcome should be
- High uncertainty indicates a knowledge gap in the model
- Reducible with more training data

# Animation: What is it?

↓ Raised on cartoons



That's animation!



That's animation!



I don't know  
what that is!



This is **Epistemic Uncertainty**.

## Lack of Knowledge:

- The man cannot tell if a stop-motion film is animation.

## Knowledge Gap:

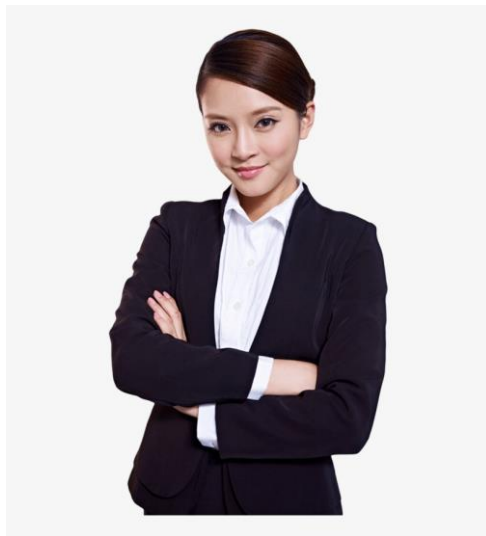
- The man was only raised (trained) on cartoon animation.
- The man could only identify cartoons as animation.

If the man had been raised (trained) on more types of animation (a broader dataset), he could've identified Coraline as animation.



# Animation: What is it?

↓ Raised on a broad variety of animation



That's animation!

*Spongebob (Cartoon Animation)*



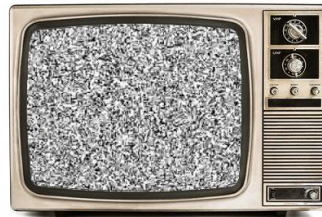
That's animation!

*Coraline (Stop-motion Animation)*



I don't know what that is!

*Someone knocked over the satellite dish*



This is **Aleatoric Uncertainty**.

## Variability/Randomness

- The satellite dish being knocked over is unrelated to the type of media being presented.

## Meaningless/Random Noise:

- The TV static does not indicate the animation type.

*Even though* the woman has been raised (trained) on more types of animation (a broader dataset), she still cannot tell if the TV program was animation or not.

# Aleatoric Uncertainty: Variance of the Input

$\mathbf{X}, \mathbf{Y}$  = Input, Target We assume  $\mathbf{Y}$  is conditionally independent given  $\mathbf{X}$

(Common assumption that  $\mathbf{Y}$  is normally distributed given  $\mathbf{X}$ .)

$P(\mathbf{Y} | \mathbf{X}) = \mathcal{N}(\mu(\mathbf{X}), \sigma^2(\mathbf{X}))$   $\mu$  and  $\sigma^2$  are the true mean and variance function

Equivalently...

$\mathbf{Y} = \mu(\mathbf{X}) + \epsilon(\mathbf{X})$  with  $\epsilon(\mathbf{X}) \sim \mathcal{N}(0, \sigma^2(\mathbf{X}))$   
plus a

$\mathbf{Y}$  is generated from  $\mathbf{X}$  by  $\mu(\mathbf{X})$

zero-mean Gaussian  $\hat{\mu}, \hat{\sigma}^2$  with variance  $\sigma^2(\mathbf{X})$ . This quantifies the input-dependent (heteroscedastic) aleatoric uncertainty.

Seitzer, M., Tavakoli, A., Antić, D., & Martius, G. (2022). On the Pitfalls of Heteroscedastic Uncertainty Estimation with Probabilistic Neural Networks. ICLR 2022 - 10th International Conference on Learning Representations. <https://arxiv.org/abs/2203.09168v2>

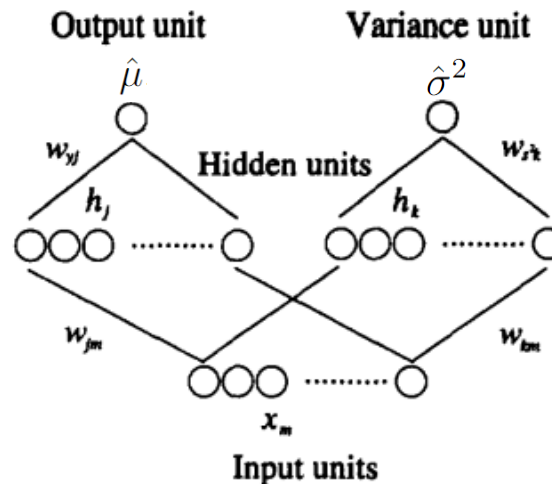
How do we get estimates ( $\mu, \sigma^2$ ) of the true mean/variance functions?

# Predicting Aleatoric Uncertainty

The NN outputs *two* values:

- predicted mean  $\hat{\mu}(\mathbf{X})$
- predicted variance  $\hat{\sigma}^2(\mathbf{X}) > 0$ .

These observed values are **treated as a sample from a Gaussian distribution** with the predicted mean and variance.



**The new predicted value (variance) must be accounted for in the loss, so that the optimizer (i.e. stochastic gradient descent) can optimize for it.**



# Predicting Aleatoric Uncertainty

Assuming the errors are normally distributed around the target:

$$p(y|x) = \frac{1}{\sqrt{2\pi\hat{\sigma}^2(x)}} e^{-\frac{(y-\hat{\mu}(x))^2}{2\hat{\sigma}^2(x)}}$$

Take the natural log of both sides to get the log-likelihood we want to maximize:

$$\ln p(y|x) = -\frac{1}{2} \ln(2\pi) - \frac{\log \hat{\sigma}^2(x)}{2} - \frac{(y - \hat{\mu}(x))^2}{2\hat{\sigma}^2(x)}$$

Negate both sides to produce the negative log-likelihood we want to minimize:

$$\mathcal{L}_{\text{NLL}} = -\ln p(y|x) = \frac{\log \hat{\sigma}^2(x)}{2} + \frac{(y - \hat{\mu}(x))^2}{2\hat{\sigma}^2(x)} + \text{constant}$$

Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2016). Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. *Advances in Neural Information Processing Systems, 2017-December*, 6403–6414. <https://arxiv.org/abs/1612.01474v3>

Seitzer, M., Tavakoli, A., Antić, D., & Martius, G. (2022). On the Pitfalls of Heteroscedastic Uncertainty Estimation with Probabilistic Neural Networks. *ICLR 2022 - 10th International Conference on Learning Representations*. <https://arxiv.org/abs/2203.09168v2>

# Predicting Aleatoric Uncertainty

The optimal parameters of the model  $\theta$  can be found using maximum likelihood estimation by minimizing the negative log-likelihood criterion  $\mathcal{L}_{\text{NLL}}$ :

$$\arg_o \min \mathcal{L}_{\text{NLL}}(\theta) = \arg_o \min_{x,y} \mathbb{E} \left[ \frac{\log \hat{\sigma}^2(x)}{2} + \frac{(y - \hat{\mu}(x))^2}{2\hat{\sigma}^2(x)} + \text{constant} \right]$$

The gradients of  $\mathcal{L}_{\text{NLL}}$  with respect to  $\hat{\mu}(x)$ ,  $\hat{\sigma}^2(x)$  are given by:

$$\nabla_{\hat{\mu}} \mathcal{L}_{\text{NLL}}(\theta) = \mathbb{E}_{x,y} \left[ \frac{\hat{\mu}(x) - y}{\hat{\sigma}^2(x)} \right], \quad \nabla_{\hat{\sigma}^2} \mathcal{L}_{\text{NLL}}(\theta) = \mathbb{E}_{x,y} \left[ \frac{\hat{\sigma}^2(x) - (y - \hat{\mu}(x))^2}{2(\hat{\sigma}^2(x))^2} \right]$$

Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2016). Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. *Advances in Neural Information Processing Systems, 2017-December*, 6403–6414. <https://arxiv.org/abs/1612.01474v3>

Seitzer, M., Tavakoli, A., Antić, D., & Martius, G. (2022). On the Pitfalls of Heteroscedastic Uncertainty Estimation with Probabilistic Neural Networks. *ICLR 2022 - 10th International Conference on Learning Representations*. <https://arxiv.org/abs/2203.09168v2>

# Estimating Aleatoric Uncertainty (Semantic Seg.)

$$\mathbf{Y} = \mu(\mathbf{X}) + \epsilon(\mathbf{X}) \quad \text{with } \epsilon(\mathbf{X}) \sim \mathcal{N}(0, \sigma^2(\mathbf{X}))$$

The objective can be approximated through Monte Carlo integration:

$$\hat{x}_{i,t} = \mu_i + \sigma_i^2 \epsilon_t \quad \text{with } \epsilon_t \sim \mathcal{N}(0, \sigma^2(x))$$

$$\mathcal{L}_x = \sum_i \log \frac{1}{T} \sum_t \exp \left( \hat{x}_{i,t,c} - \log \sum_{c'} \exp \hat{x}_{i,t,c'} \right)$$

In the following code example, we use cross-entropy loss (which implicitly applies a softmax activation) to create a semantic segmentation aleatoric uncertainty-aware loss function.

```
# Fetch training data and pass it through the model
T = 2
criterion = torch.nn.CrossEntropyLoss()
data, target = ...
mu, sigma_squared = model(data)

# Obtain the loss
p_hat = []
for t in range(T):
    epsilon = torch.normal(mean=torch.zeros(mu.shape),
                           std=torch.ones(mu.shape))
    noisy_output = mu + sigma_squared * epsilon
    p_hat.append(criterion(noisy_output, target))
loss = torch.mean(torch.stack(p_hat, dim=-1))
loss.backward()
```

# Epistemic Uncertainty: Uncertainty of the Model

Model uncertainty indicates a knowledge gap in the model.

By repeatedly sampling a model with different dropout parameters, this uncertainty can be represented.

## Monte Carlo (MC) Dropout

- Does not require a “special” loss function
- Can be used with existing NN models trained with dropout
- The multiple passes can (with good enough hardware) be done concurrently

# Epistemic Uncertainty Using MC Dropout

With  $\mathcal{W} = \{\mathbf{W}_i\}_{i=1}^L$  as the set of dropout-enabled weights for a model with  $L$  layers and  $q(\mathcal{W})$  defined as an approximating variational distribution:

$$q(\mathbf{y}^*|x^*) = \int p(\mathbf{y}^*|x^*, \mathcal{W}) q(\mathcal{W}) d\mathcal{W}$$

The estimate can be acquired through a sampling of  $T$  groups of dropout-enabled weights:

$$\mathbb{E}(\mathbf{y}^*) \approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}^*(x^*, \mathbf{W}_1^t, \dots, \mathbf{W}_L^t)$$

This is equivalent to performing  $T$  stochastic forward passes through the network with dropouts enabled and averaging the results.

The sample variance across  $T$  stochastic forward passes through the network with dropouts enabled can also be calculated to determine the model uncertainty.



- Seitzer, M., Tavakoli, A., Antić, D., & Martius, G. (2022). On the Pitfalls of Heteroscedastic Uncertainty Estimation with Probabilistic Neural Networks. *ICLR 2022 - 10th International Conference on Learning Representations*. <https://arxiv.org/abs/2203.09168v2>
- Nix, D. A., & Weigend, A. S. (1994). Estimating the mean and variance of the target probability distribution. *IEEE International Conference on Neural Networks - Conference Proceedings, 1*, 55–60. <https://doi.org/10.1109/ICNN.1994.374138>
- Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2016). Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. *Advances in Neural Information Processing Systems, 2017-December*, 6403–6414. <https://arxiv.org/abs/1612.01474v3>
- Kendall, A., & Gal, Y. (2017). What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? *Advances in Neural Information Processing Systems, 2017-December*, 5575–5585. <https://arxiv.org/abs/1703.04977v2>
- Gal, Y., & Ghahramani, Z. (2015). Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. *33rd International Conference on Machine Learning, ICML 2016, 3*, 1651–1660. <https://arxiv.org/abs/1506.02142v6>