

What my deep model doesn't know...

Yarin Gal

yg279@cam.ac.uk

Why should I care about uncertainty?

- ▶ We train a model to recognise dog breeds



Why should I care about uncertainty?

- ▶ We train a model to recognise dog breeds
- ▶ And are given a cat to classify



- ▶ We train a model to recognise dog breeds
- ▶ And are given a cat to classify
- ▶ What would you want your model to do?



- ▶ We train a model to recognise dog breeds
- ▶ And are given a cat to classify
- ▶ What would you want your model to do?
- ▶ Similar problems in *decision making, physics, life science, etc.*¹



- ▶ We train a model to recognise dog breeds
- ▶ And are given a cat to classify
- ▶ What would you want your model to do?
- ▶ Similar problems in *decision making, physics, life science, etc.*¹
- ▶ For the practitioner: model debugging, specialised models, critical systems



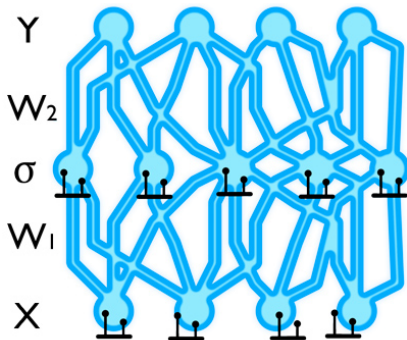
- ▶ We train a model to recognise dog breeds
- ▶ And are given a cat to classify
- ▶ What would you want your model to do?
- ▶ Similar problems in *decision making, physics, life science, etc.*¹
- ▶ For the practitioner: model debugging, specialised models, critical systems
- ▶ We need a way to tell **what our model knows** and what not.
- ▶ Luckily, if you use **dropout** you already have this information. You just **ignored** it so far.

¹Complete references at end of slides

- ▶ We train a model to recognise dog breeds
- ▶ And are given a cat to classify
- ▶ What would you want your model to do?
- ▶ Similar problems in *decision making, physics, life science, etc.*¹
- ▶ For the practitioner: model debugging, specialised models, critical systems
- ▶ We need a way to tell **what our model knows** and what not.
- ▶ Luckily, if you use **dropout** you already have this information. You just **ignored** it so far.

¹Complete references at end of slides

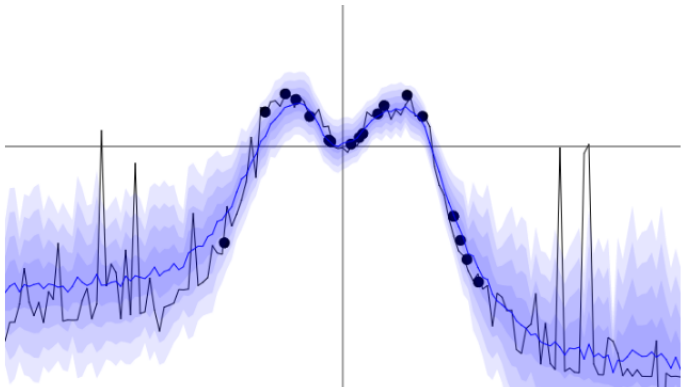
- ▶ Used in **most modern deep learning models**
- ▶ It circumvents **over-fitting**
- ▶ And improves **performance**



- ▶ **Training time:** drop units, **test time:** don't drop

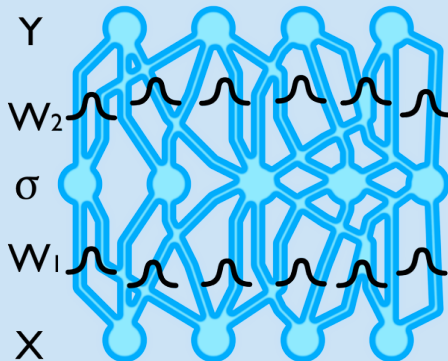
Bayesian modelling:

- ▶ Observed inputs $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$ and outputs $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^N$
- ▶ Capture distribution believed to have generated outputs
- ▶ Look at the first two moments:



From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s



- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$

From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...

- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$

- ▶ **Inference** with

$$q_{\theta}(\mathbf{W}) := \mathbf{M} \cdot \text{diag}(\text{Bernoulli})$$

and parameter \mathbf{M}

= **Dropout training.**

From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...

- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$

- ▶ **Inference** with

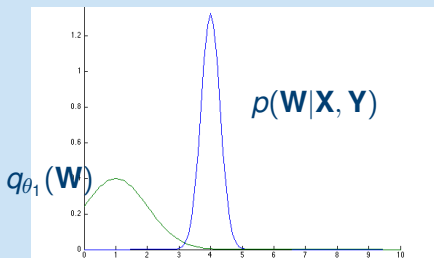
$$q_{\theta}(\mathbf{W}) := \mathbf{M} \cdot \text{diag}(\text{Bernoulli})$$

and parameter \mathbf{M}

= **Dropout training.**

From Bayesian neural networks to Dropout

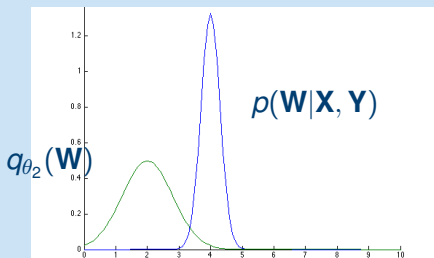
- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...
- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate
$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$



From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...
- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

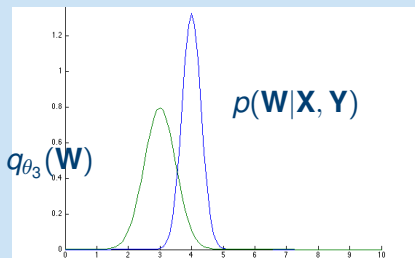
$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$



From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...
- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

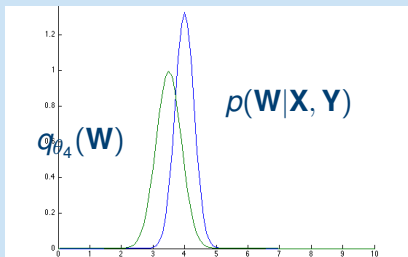
$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$



From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...
- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

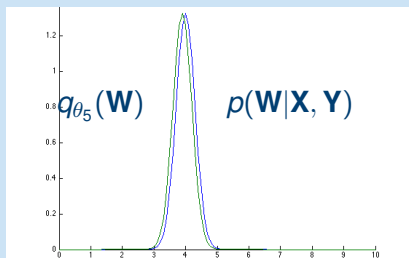
$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$



From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...
- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$



From Bayesian neural networks to Dropout

- ▶ Place **prior** $p(\mathbf{W})$ dist. on weights, making these r.v.s
- ▶ Given dataset \mathbf{X}, \mathbf{Y} , the r.v. \mathbf{W} has a **posterior**: $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$
- ▶ Which is difficult to evaluate...

- ▶ Can define **simple distribution** $q_{\theta}(\cdot)$ and approximate

$$q_{\theta}(\mathbf{W}) \approx p(\mathbf{W}|\mathbf{X}, \mathbf{Y}).$$

- ▶ **Inference** with

$$q_{\theta}(\mathbf{W}) := \mathbf{M} \cdot \text{diag}(\text{Bernoulli})$$

and parameter \mathbf{M}

= **Dropout training.**

The theory above means that with dropout we:

- ▶ **capture distribution** that generated observed data
- ▶ can combine model with Bayesian techniques in a **practical** way...
- ▶ have **uncertainty estimates** in the network

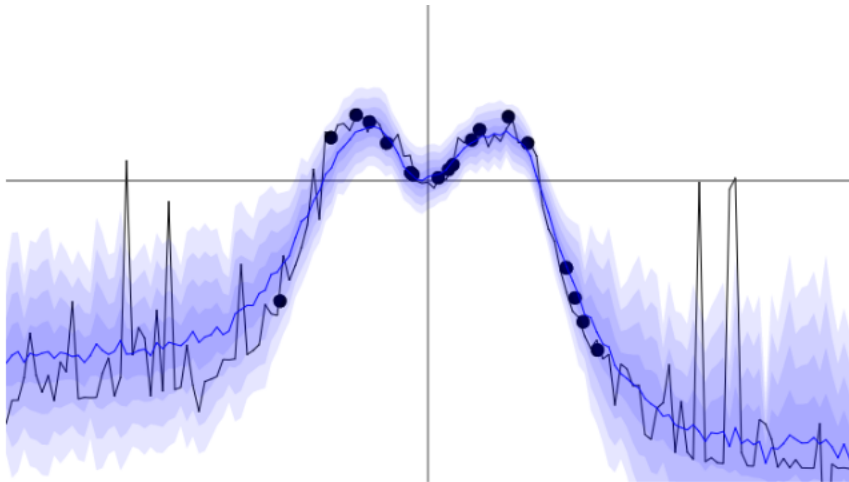
The theory above means that with dropout we:

- ▶ **capture distribution** that generated observed data
- ▶ can combine model with Bayesian techniques in a **practical** way...
- ▶ have **uncertainty estimates** in the network

The theory above means that with dropout we:

- ▶ **capture distribution** that generated observed data
- ▶ can combine model with Bayesian techniques in a **practical** way...
- ▶ have **uncertainty estimates** in the network

We fit a **distribution**...



We fit a **distribution**...

- ▶ Use first moment for **predictions**:

$$\mathbb{E}(\mathbf{y}^*) \approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}_t$$

with $\hat{\mathbf{y}}_t \sim \text{DropoutNetwork}(\mathbf{x}^*)$.

- ▶ Use second moment for **uncertainty** (in regression):

$$\text{Var}(\mathbf{y}^*) \approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}_t^T \hat{\mathbf{y}}_t - \mathbb{E}(\mathbf{y}^*)^T \mathbb{E}(\mathbf{y}^*) + \tau^{-1} \mathbf{I}$$

with $\hat{\mathbf{y}}_t \sim \text{DropoutNetwork}(\mathbf{x}^*)$.

We fit a **distribution**...

- ▶ Use first moment for **predictions**:

$$\mathbb{E}(\mathbf{y}^*) \approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}_t$$

with $\hat{\mathbf{y}}_t \sim \text{DropoutNetwork}(\mathbf{x}^*)$.

- ▶ Use second moment for **uncertainty** (in regression):

$$\text{Var}(\mathbf{y}^*) \approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}_t^T \hat{\mathbf{y}}_t - \mathbb{E}(\mathbf{y}^*)^T \mathbb{E}(\mathbf{y}^*) + \tau^{-1} \mathbf{I}$$

with $\hat{\mathbf{y}}_t \sim \text{DropoutNetwork}(\mathbf{x}^*)$.

In more practical terms, given point x :²

- ▶ drop units at test time
- ▶ repeat 10 times
- ▶ and look at mean and sample variance.
- ▶ Or in Python:

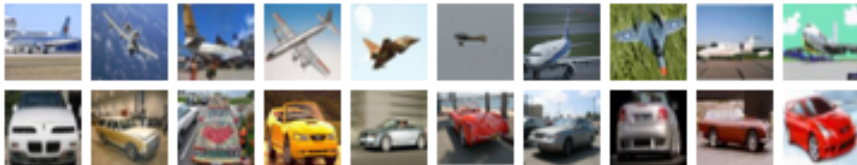
```
1 | y = []
2 | for _ in xrange(10):
3 |     y.append(model.output(x, dropout=True))
4 | y_mean = numpy.mean(y)
5 | y_var = numpy.var(y)
```

²Friendly introduction given in yarin.co/blog

CIFAR Test Error (and Std.)

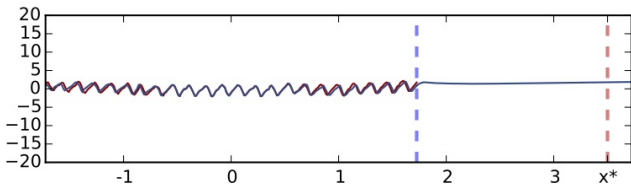
Model	Standard Dropout	Bayesian technique
NIN	10.43 (Lin et al., 2013)	10.27 \pm 0.05
DSN	9.37 (Lee et al., 2014)	9.32 \pm 0.02
Augmented-DSN	7.95 (Lee et al., 2014)	7.71 \pm 0.09

Table : Bayesian techniques with existing state-of-the-art



What would be the CO_2 concentration level in Mauna Loa, Hawaii, *in 20 years' time*?

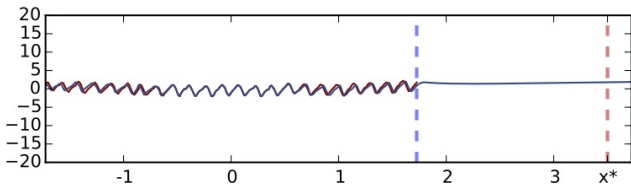
- ▶ Normal dropout (weight averaging, 5 layers, ReLU units):



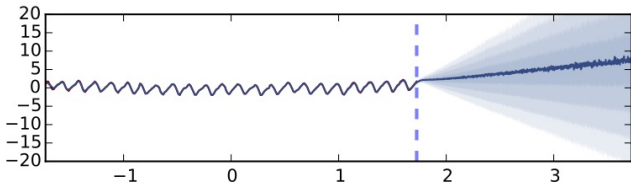
- ▶ Same network, Bayesian perspective:

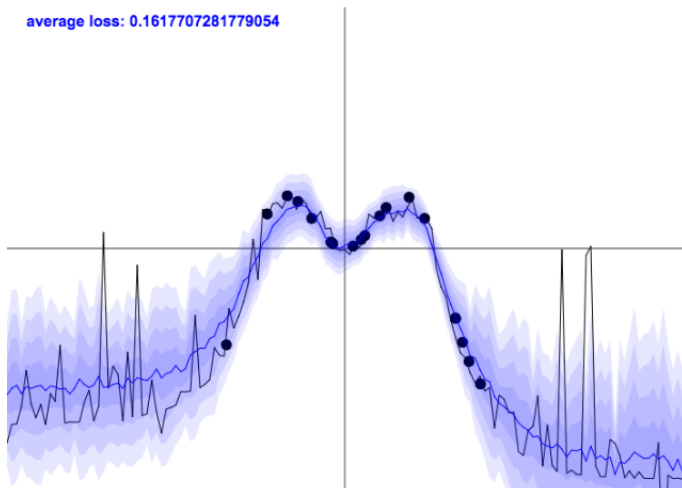
What would be the CO_2 concentration level in Mauna Loa, Hawaii, *in 20 years' time*?

- ▶ Normal dropout (weight averaging, 5 layers, ReLU units):



- ▶ Same network, Bayesian perspective:



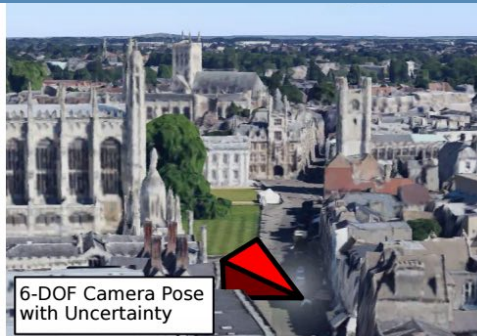
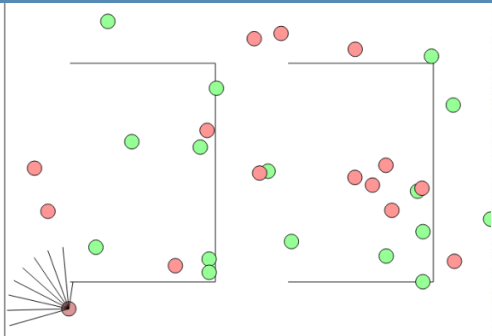


[Online demo]³

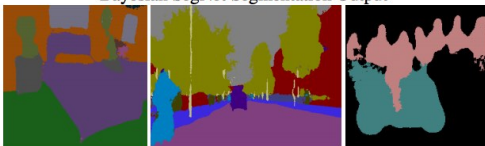
Dataset	Avg. Test RMSE and Std. Errors			Avg. Test LL and Std. Errors		
	VI	PBP	Dropout	VI	PBP	Dropout
Boston Housing	4.32 \pm 0.29	3.01 \pm 0.18	2.97 \pm0.85	-2.90 \pm 0.07	-2.57 \pm 0.09	-2.46 \pm0.25
Concrete Strength	7.19 \pm 0.12	5.67 \pm 0.09	5.23 \pm0.53	-3.39 \pm 0.02	-3.16 \pm 0.02	-3.04 \pm0.09
Energy Efficiency	2.65 \pm 0.08	1.80 \pm 0.05	1.66 \pm0.19	-2.39 \pm 0.03	-2.04 \pm 0.02	-1.99 \pm0.09
Kin8nm	0.10 \pm0.00	0.10 \pm0.00	0.10 \pm0.00	0.90 \pm 0.01	0.90 \pm 0.01	0.95 \pm0.03
Naval Propulsion	0.01 \pm0.00	0.01 \pm0.00	0.01 \pm0.00	3.73 \pm 0.12	3.73 \pm 0.01	3.80 \pm0.05
Power Plant	4.33 \pm 0.04	4.12 \pm 0.03	4.02 \pm0.18	-2.89 \pm 0.01	-2.84 \pm 0.01	-2.80 \pm0.05
Protein Structure	4.84 \pm 0.03	4.73 \pm 0.01	4.36 \pm0.04	-2.99 \pm 0.01	-2.97 \pm 0.00	-2.89 \pm0.01
Wine Quality Red	0.65 \pm 0.01	0.64 \pm 0.01	0.62 \pm0.04	-0.98 \pm 0.01	-0.97 \pm 0.01	-0.93 \pm0.06
Yacht Hydrodynamics	6.89 \pm 0.67	1.02 \pm0.05	1.11 \pm 0.38	-3.43 \pm 0.16	-1.63 \pm 0.02	-1.55 \pm0.12
Year Prediction MSD	9.034 \pm NA	8.879 \pm NA	8.849 \pmNA	-3.622 \pm NA	-3.603 \pm NA	-3.588 \pmNA

Table 1: **Average test performance in RMSE and predictive log likelihood** for a popular variational inference method (VI, Graves [20]), Probabilistic back-propagation (PBP, Hernández-Lobato and Adams [27]), and dropout uncertainty (Dropout).

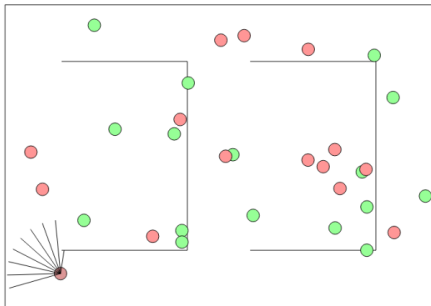
Applications



Bayesian SegNet Segmentation Output



- ▶ We have a “Roomba”⁴
- ▶ Penalised -5 for walking into a wall, $+10$ reward for collecting dirt
- ▶ Our environment is stochastic and ever changing
- ▶ We want a net to learn what actions to do in different situations



⁴Code based on Karpathy and authors. github.com/karpathy/convnetjs

Behavioural policies:

- ▶ **Epsilon-greedy** – take random actions with probability ϵ and optimal actions otherwise
- ▶ Using uncertainty we can learn faster
- ▶ **Thompson sampling** – draw realisation from current belief over world, choose action with highest value
- ▶ In practice: simulate a stochastic forward pass through the dropout network and choose action with highest value

Behavioural policies:

- ▶ **Epsilon-greedy** – take random actions with probability ϵ and optimal actions otherwise
- ▶ Using uncertainty we can learn faster
- ▶ **Thompson sampling** – draw realisation from current belief over world, choose action with highest value
- ▶ In practice: simulate a stochastic forward pass through the dropout network and choose action with highest value

Behavioural policies:

- ▶ **Epsilon-greedy** – take random actions with probability ϵ and optimal actions otherwise
- ▶ Using uncertainty we can learn faster
- ▶ **Thompson sampling** – draw realisation from current belief over world, choose action with highest value
- ▶ In practice: simulate a stochastic forward pass through the dropout network and choose action with highest value

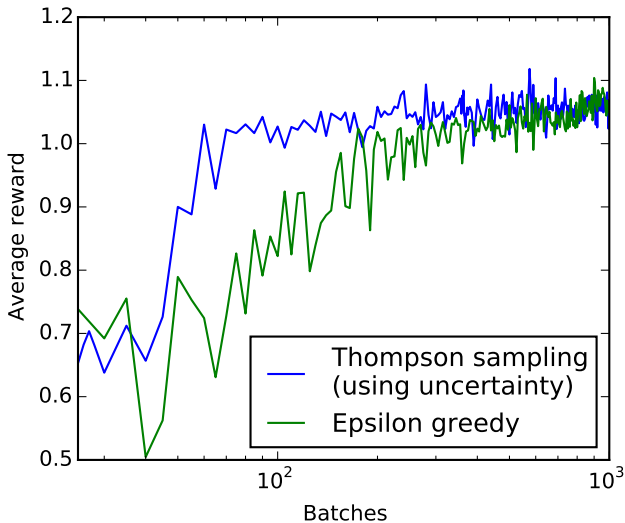
Behavioural policies:

- ▶ **Epsilon-greedy** – take random actions with probability ϵ and optimal actions otherwise
- ▶ Using uncertainty we can learn faster
- ▶ **Thompson sampling** – draw realisation from current belief over world, choose action with highest value
- ▶ In practice: simulate a stochastic forward pass through the dropout network and choose action with highest value



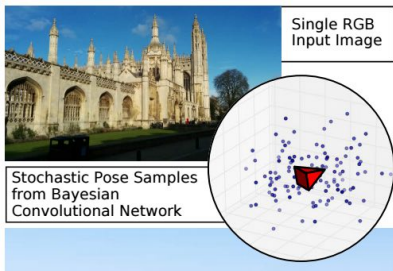
[\[Online demo\]](#)⁵

⁵yarin.co/blog



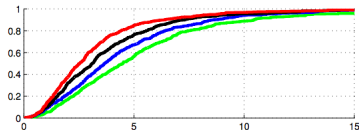
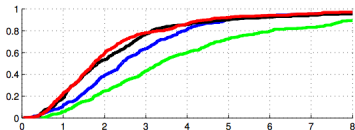
Average reward over time (log scale)

- ▶ Where was a picture taken? (Kendall and Cipolla, 2015)⁶



- ▶ With Bayesian techniques above: **10–15%** improvement on **state-of-the-art**
- ▶ Uncertainty increases as a test photo diverges from training distribution
- ▶ Test photos with high uncertainty (strong occlusion from vehicles, pedestrians or other objects)

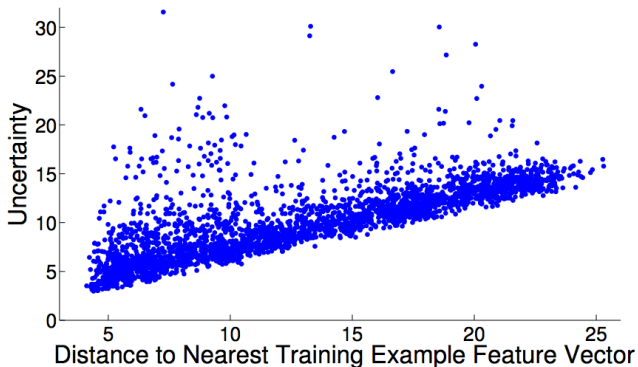
- ▶ Where was a picture taken? (Kendall and Cipolla, 2015)⁶
- ▶ With Bayesian techniques above: **10–15% improvement on state-of-the-art**



- ▶ Uncertainty increases as a test photo diverges from training distribution
- ▶ Test photos with high uncertainty (strong occlusion from vehicles, pedestrians or other objects)
- ▶ Localisation error correlates with uncertainty

⁶Figures used with author permission

- ▶ Where was a picture taken? (Kendall and Cipolla, 2015)⁶
- ▶ With Bayesian techniques above: **10–15% improvement on state-of-the-art**
- ▶ Uncertainty increases as a test photo diverges from training distribution

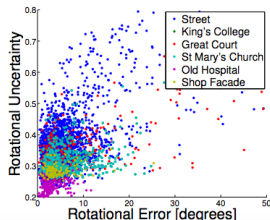
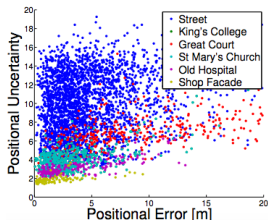


- ▶ Where was a picture taken? (Kendall and Cipolla, 2015)⁶
- ▶ With Bayesian techniques above: **10–15% improvement on state-of-the-art**
- ▶ Uncertainty increases as a test photo diverges from training distribution
- ▶ Test photos with high uncertainty (strong occlusion from vehicles, pedestrians or other objects)



- ▶ Localisation error correlates with uncertainty

- ▶ Where was a picture taken? (Kendall and Cipolla, 2015)⁶
- ▶ With Bayesian techniques above: **10–15% improvement on state-of-the-art**
- ▶ Uncertainty increases as a test photo diverges from training distribution
- ▶ Test photos with high uncertainty (strong occlusion from vehicles, pedestrians or other objects)
- ▶ Localisation error correlates with uncertainty

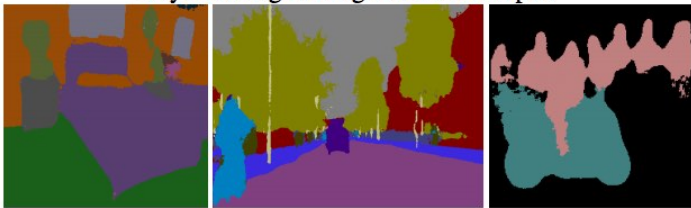


- Scene understanding: what's in a photo and where? (Kendall, Badrinarayanan, and Cipolla, 2015)⁷

Input Images

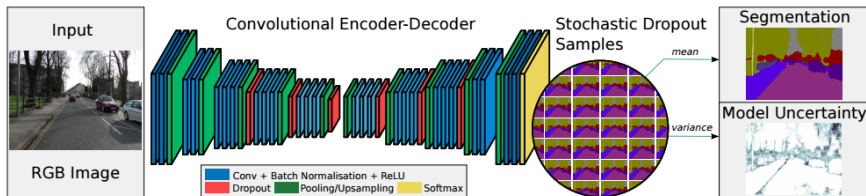


Bayesian SegNet Segmentation Output



⁷Figures used with author permission

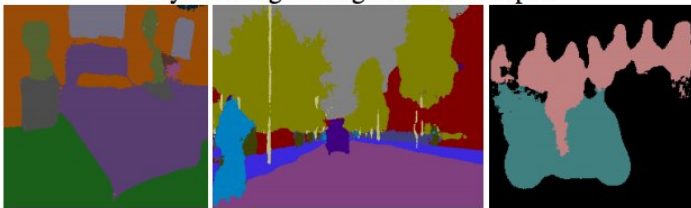
- Scene understanding: what's in a photo and where? (Kendall, Badrinarayanan, and Cipolla, 2015)⁷



⁷Figures used with author permission

- Scene understanding: what's in a photo and where? (Kendall, Badrinarayanan, and Cipolla, 2015)⁷

Bayesian SegNet Segmentation Output

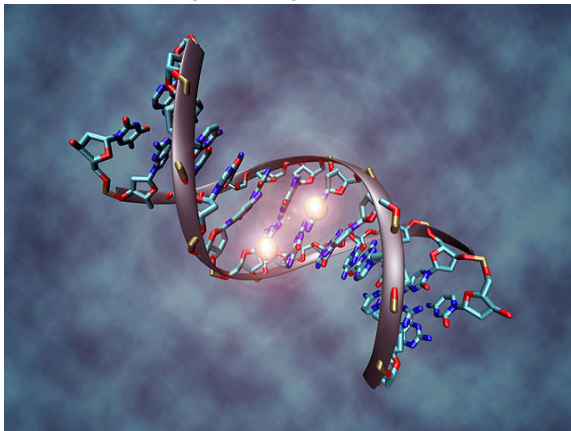


Bayesian SegNet Model Uncertainty Output



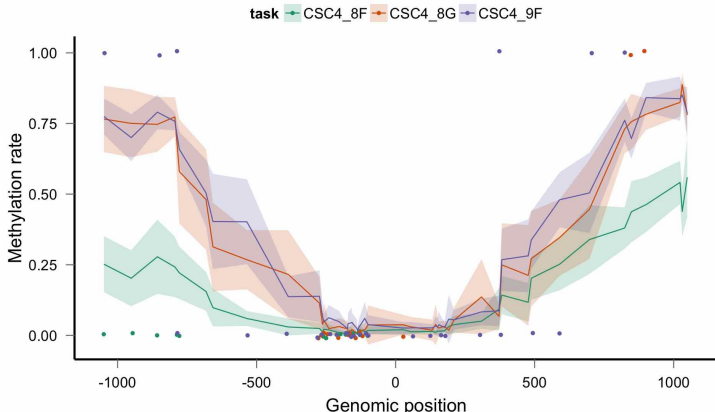
⁷Figures used with author permission

- ▶ Angermueller and Stegle (2015) fit a network to predict **DNA methylation** – controls gene regulation



- ▶ Look at methylation rate of different embryonic stem cells. **Uncertainty increases** in genomic contexts that are hard to predict (e.g. LMR or H3K27me3)

- ▶ Angermueller and Stegle (2015) fit a network to predict **DNA methylation** – controls gene regulation
- ▶ Look at methylation rate of different embryonic stem cells. **Uncertainty increases** in genomic contexts that are hard to predict (e.g. LMR or H3K27me3)



If you use **dropout** you already have uncertainty information = **practical** deep learning uncertainty.

- ▶ Applications: capture language ambiguity?

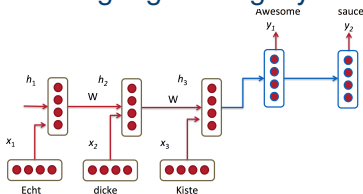
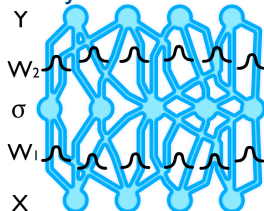


Image Source: cs224d.stanford.edu/lectures/CS224d-Lecture8.pdf

- ▶ Tools: weight uncertainty for model debugging?



If you use **dropout** you already have uncertainty information = **practical** deep learning uncertainty.

- ▶ Applications: capture language ambiguity?

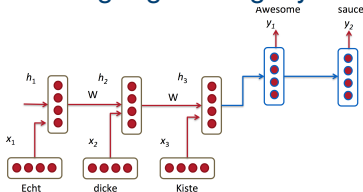
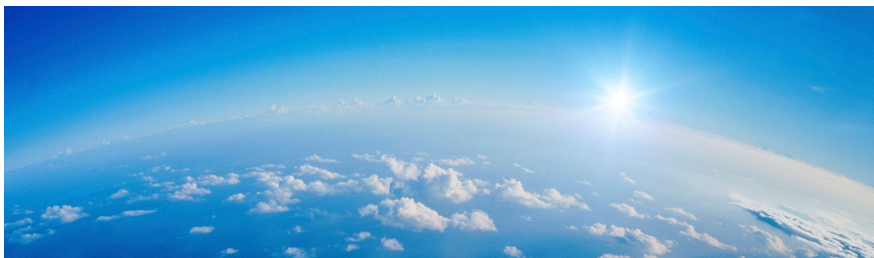


Image Source: cs224d.stanford.edu/lectures/CS224d-Lecture8.pdf

- ▶ Tools: weight uncertainty for model debugging?

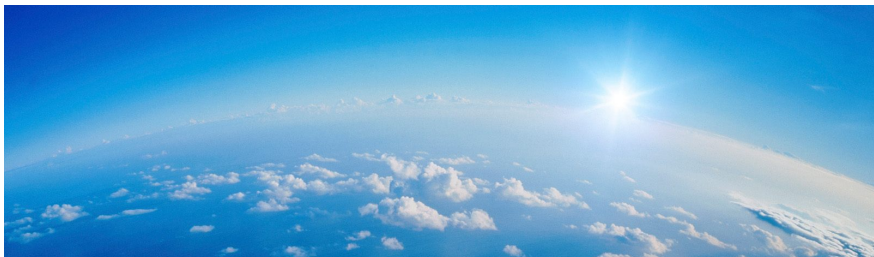
Work in progress!



Most exciting is work to come:

- ▶ **Practical uncertainty** in deep learning applications
- ▶ **Principled extensions** to deep learning tools
- ▶ **Hybrid** deep learning – Bayesian models

and much, much, more.



Most exciting is work to come:

- ▶ **Practical uncertainty** in deep learning applications
- ▶ **Principled extensions** to deep learning tools
- ▶ **Hybrid** deep learning – Bayesian models

and much, much, more.

Thank you for listening.

- ▶ Y Gal, Z Ghahramani, “**Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning**”, arXiv preprint, arXiv:1506.02142 (2015).
- ▶ Y Gal, Z Ghahramani, “**Dropout as a Bayesian Approximation: Appendix**”, arXiv preprint, arXiv:1506.02157 (2015).
- ▶ Y Gal, Z Ghahramani, “**Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference**”, arXiv preprint, arXiv:1506.02158 (2015).
- ▶ A Kendall, R Cipolla, “**Modelling Uncertainty in Deep Learning for Camera Relocalization**”, arXiv preprint, arXiv:1509.05909 (2015)
- ▶ C Angermueller and O Stegle, “**Multi-task deep neural network to predict CpG methylation profiles from low-coverage sequencing data**”, NIPS MLCB workshop (2015).
- ▶ JM Hernandez-Lobato, RP Adams, “**Probabilistic Backpropagation for Scalable Learning of Bayesian Neural Networks**”, ICML (2015).
- ▶ DP Kingma, T Salimans, M Welling, “**Variational Dropout and the Local Reparameterization Trick**”, NIPS (2015).
- ▶ DJ Rezende, S Mohamed, D Wierstra, “**Stochastic Backpropagation and Approximate Inference in Deep Generative Models**”, ICML (2014).

- ▶ Krzywinski and Altman, “**Points of significance: Importance of being uncertain**”, Nature Methods (2013).
- ▶ Herzog and Ostwald, “**Experimental biology: Sometimes Bayesian statistics are better**”, Nature (2013).
- ▶ Nuzzo, “**Scientific method: Statistical errors**”, Nature (2014).
- ▶ Woolston, “**Psychology journal bans P values**”, Nature (2015).
- ▶ Ghahramani, “**Probabilistic machine learning and artificial intelligence**”, Nature (2015).