# Dropout as a Structured Shrinkage Prior

## Eric Nalisnick



UNIVERSITY OF

CAMBRIDGE



Computational and Biological Learning University of Cambridge

Randomly set hidden units to zero (i.e. drop them) for every forward pass during training [Hinton et al., 2012].



Standard Neural Network



After Applying Dropout

Randomly set hidden units to zero (i.e. drop them) for every forward pass during training [Hinton et al., 2012].



Randomly set hidden units to zero (i.e. drop them) for every forward pass during training [Hinton et al., 2012].



Randomly set hidden units to zero (i.e. drop them) for every forward pass during training [Hinton et al., 2012].





Geoffrey Hinton 2018 Turing Award winner for Deep Learning

Source: https://www.reddit.com/r/MachineLearning/comments/4w6tsv/ama\_we\_are\_the\_google\_brain\_team\_wed\_love\_to/d6dgyse

I went to my bank. The tellers kept changing.... figured it must be because it would require cooperation between employees to successfully defraud the bank. This made me realize that randomly removing... neurons...would prevent conspiracies and thus reduce overfitting.



I went to my bank. The tellers kept changing....l figured it must be because it would require cooperation between employees to successfully defraud the bank. This made me realize that randomly removing... neurons...would prevent conspiracies and thus reduce overfitting.



I went to my bank. The tellers kept changing....I figured it must be because it would require cooperation between employees to successfully defraud the bank. This made me realize that randomly removing... neurons...would prevent conspiracies and thus reduce overfitting.



I went to my bank. The tellers kept changing....I figured it must be because it would require cooperation between employees to successfully defraud the bank. This made me realize that randomly removing neurons...would prevent conspiracies and thus reduce overfitting.



## $\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\mathbf{W}_l)$



Implementation as *multiplicative noise*:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\mathbf{\Lambda}_l\mathbf{W}_l)$$

Implementation as *multiplicative noise*:

 $\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\boldsymbol{\Lambda}_l \mathbf{W}_l)$ 

 $\lambda_{l.i.i} \sim p(\lambda)$ 

Diagonal Matrix of Random Variables

Implementation as *multiplicative noise*:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\boldsymbol{\Lambda}_l \mathbf{W}_l)$$

 $\lambda_{l,i,i} \sim p(\lambda)$ 

Diagonal Matrix of Random Variables

 $\otimes$  Dropout corresponds to  $p(\lambda)$  being Bernoulli.

Implementation as *multiplicative noise*:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\boldsymbol{\Lambda}_l \mathbf{W}_l)$$

$$\lambda_{l,i,i} \sim p(\lambda)$$

Diagonal Matrix of Random Variables

- $\otimes$  Dropout corresponds to  $p(\lambda)$  being Bernoulli.
- $\otimes$  Gaussian, beta, and uniform noise also work well.

Implementation as *multiplicative noise*:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\boldsymbol{\Lambda}_l \mathbf{W}_l)$$

$$\lambda_{l,i,i} \sim p(\lambda)$$

Diagonal Matrix of Random Variables

- $\otimes$  Dropout corresponds to  $p(\lambda)$  being Bernoulli.
- $\otimes$  Gaussian, beta, and uniform noise also work well.
- Optimization objective:

$$\mathbb{E}_{\lambda}\left[\log p\left(\mathbf{y} | \mathbf{X}, \left\{\mathbf{W}_{l}\right\}_{l=1}^{L+1}, \left\{\mathbf{\Lambda}_{l}\right\}_{l=1}^{L+1}\right)\right]$$

Solution Series Seri

 $p(\lambda)$  serves as a variational approximation to the posterior  $p(\lambda | \mathbf{y}, \mathbf{X})$ 

- ⊗ Gal & Ghahramani [2016] argued that dropout can be interpreted as variational Bayesian inference.
- Obtain predictive uncertainty by averaging over noise samples.

$$p(y^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) \approx \frac{1}{S} \sum_{s=1}^{S} p(y^* | \mathbf{x}^*, \{\hat{\boldsymbol{\Lambda}}_{l,s}\}_{l=1}^{L+1}), \quad \hat{\lambda} \sim p(\lambda)$$

- Solution Series Seri
- Obtain predictive uncertainty by averaging over noise samples.

$$p(y^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) \approx \frac{1}{S} \sum_{s=1}^{S} p(y^* | \mathbf{x}^*, \{\hat{\mathbf{\Lambda}}_{l,s}\}_{l=1}^{L+1}), \quad \hat{\lambda} \sim p(\lambda)$$



- Solution Series Seri
- Obtain predictive uncertainty by averaging over noise samples.
- Solution  $\otimes$  However, p(λ)—their posterior approximation—is fixed, does not depend on data.

- Solution Series Seri
- Obtain predictive uncertainty by averaging over noise samples.
- $\otimes$  However, p(λ)—their posterior approximation—is fixed, does not depend on data.
- Subsequent work has attempted to fix this by optimizing the noise distribution. [Maeda, 2014; Kingma et al., 2015; Molchanov et al., 2017; Gal et al., 2017]

THE ASTROPHYSICAL JOURNAL, 859:64 (16pp), 2018 May 20 © 2018. The American Astronomical Society. All rights reserved. https://doi.org/10.3847/1538-4357/aabfdb



#### Detecting Solar-like Oscillations in Red Giants with Deep Learning

Marc Hon<sup>1</sup><sup>(1)</sup>, Dennis Stello<sup>1,2,3</sup><sup>(1)</sup>, and Joel C. Zinn<sup>4</sup><sup>(1)</sup>

<sup>1</sup> School of Physics, The University of New South Wales, Sydney NSW 2052, Australia; mtyh555@uowmail.edu.au
 <sup>2</sup> Sydney Institute for Astronomy (SIfA), School of Physics, University of Sydney, NSW 2006, Australia
 <sup>3</sup> Stellar Astrophysics Centre, Department of Physics and Astronomy, Aarhus University, Ny Munkegade 120, DK-8000 Aarhus C, Denmark
 <sup>4</sup> Department of Astronomy, The Ohio State University, Columbus, OH 43210, USA
 *Received 2018 March 12; revised 2018 April 18; accepted 2018 April 19; published 2018 May 24*

THE ASTROPHYSICAL JOURNAL, 859:64 (16pp), 2018 May 20 © 2018. The American Astronomical Society. All rights reserved. https://doi.org/10.3847/1538-4357/aabfdb



#### Detecting Solar-like Oscillations in Red Giants with Deep Learning

Marc Hon<sup>1</sup><sup>(1)</sup>, Dennis Stello<sup>1,2,3</sup><sup>(1)</sup>, and Joel C. Zinn<sup>4</sup><sup>(1)</sup>

<sup>1</sup> School of Physics, The University of New South Wales, Sydney NSW 2052, Australia; mtyh555@uowmail.edu.au
 <sup>2</sup> Sydney Institute for Astronomy (SIfA), School of Physics, University of Sydney, NSW 2006, Australia
 <sup>3</sup> Stellar Astrophysics Centre, Department of Physics and Astronomy, Aarhus University, Ny Munkegade 120, DK-8000 Aarhus C, Denmark
 <sup>4</sup> Department of Astronomy, The Ohio State University, Columbus, OH 43210, USA
 *Received 2018 March 12; revised 2018 April 18; accepted 2018 April 19; published 2018 May 24*

"...[dropout] is effectively a Monte Carlo integration over a Gaussian process posterior approximation (Gal & Ghahramani, 2016)."

~ Hon et al. (2018)

Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not restricted to variational inference).

- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not restricted to variational inference).
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.

- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not restricted to variational inference).
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.
- Principles for extension to new architectures?

- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not restricted to variational inference).
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.
- ⊗ Principles for extension to new architectures?

Subscription of the dropout / mult. noise equation:  $h = f(h + \lambda w) = \lambda - n(\lambda)$ 

$$h_l = f_l(h_{l-1}\lambda_l w_l), \quad \lambda_l \sim p(\lambda)$$

- ⊗ Consider a scalar-version of the dropout / mult. noise equation:  $h_l = f_l(h_{l-1}\lambda_l w_l), \quad \lambda_l \sim p(\lambda)$
- $\otimes$  Assume a Gaussian prior on the weight:

$$w_l \sim \mathbf{N}(0, \sigma_0^2)$$

- Solution Consider a scalar-version of the dropout / mult. noise equation:  $h_l = f_l(h_{l-1}\lambda_l w_l), \quad \lambda_l \sim p(\lambda)$
- $\otimes$  Assume a Gaussian prior on the weight:

$$w_l \sim \mathbf{N}(0, \sigma_0^2)$$

What is the distribution of the noise-weight product?

$$h_l = f_l(h_{l-1}\lambda_l w_l)$$

Distribution of this product?

- Solution
  Solution
  Solution
  Solution
  Solution
  Model  $h_l = f_l(h_{l-1}\lambda_l w_l), \quad \lambda_l \sim p(\lambda)$
- $\otimes$  Assume a Gaussian prior on the weight:

$$w_l \sim \mathbf{N}(0, \sigma_0^2)$$

⊗ What is the distribution of the noise-weight product?

$$h_l = f_l(h_{l-1}\lambda_l w_l)$$

Gaussian scale mixture! (expanded param.) [Beale & Mallows, 1959]

Gaussian scale mixtures can be reparameterized into the following hierarchical form:

$$h_l = f_l(h_{l-1}\lambda_l w_l), \quad w_l \sim \mathbf{N}(0,\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

Gaussian scale mixtures can be reparameterized into the following hierarchical form:

$$h_l = f_l(h_{l-1}\lambda_l w_l), \quad w_l \sim \mathbf{N}(0,\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

SWITCH TO HIERARCHICAL PARAMETRIZATION

$$h_l = f_l(h_{l-1}\tilde{w}_l), \quad \tilde{w}_l \sim \mathbf{N}(0,\lambda_l^2\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

Gaussian scale mixtures can be reparameterized into the following hierarchical form:

$$h_l = f_l(h_{l-1}\lambda_l w_l), \quad w_l \sim \mathbf{N}(0,\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

SWITCH TO HIERARCHICAL PARAMETRIZATION

$$h_l = f_l(h_{l-1}\tilde{w}_l), \quad \tilde{w}_l \sim \mathbf{N}(0,\lambda_l^2\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

 $\otimes$  Noise moves from likelihood and becomes a scale
#### Dropout as a Bayesian Prior

Gaussian scale mixtures can be reparameterized into the following hierarchical form:

$$h_l = f_l(h_{l-1}\lambda_l w_l), \quad w_l \sim \mathbf{N}(0,\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

SWITCH TO HIERARCHICAL PARAMETRIZATION

$$h_l = f_l(h_{l-1}\tilde{w}_l), \quad \tilde{w}_l \sim \mathbf{N}(0,\lambda_l^2\sigma_0^2), \quad \lambda_l \sim p(\lambda)$$

- ⊗ Noise moves from likelihood and becomes a scale
- Second States State

#### Dropout as a Bayesian Prior

 Reparameterization translates between noise distributions and marginal priors:

Noise Model $p(\lambda)$	Variance Prior $p(\lambda^2)$	Marginal Prior $p(\widetilde{w})$
Bernoulli	Bernoulli	Spike-and-Slab
Gaussian	$\chi^2$	Generalized Hyperbolic
Rayleigh	Exponential	Laplace
Inverse Nakagami	$\Gamma^{-1}$	Student-t
Half-Cauchy	Unnamed	Horseshoe

THE ASTROPHYSICAL JOURNAL, 859:64 (16pp), 2018 May 20 © 2018. The American Astronomical Society. All rights reserved. https://doi.org/10.3847/1538-4357/aabfdb



#### Detecting Solar-like Oscillations in Red Giants with Deep Learning

Marc Hon<sup>1</sup><sup>(1)</sup>, Dennis Stello<sup>1,2,3</sup><sup>(1)</sup>, and Joel C. Zinn<sup>4</sup><sup>(1)</sup>

<sup>1</sup> School of Physics, The University of New South Wales, Sydney NSW 2052, Australia; mtyh555@uowmail.edu.au
 <sup>2</sup> Sydney Institute for Astronomy (SIfA), School of Physics, University of Sydney, NSW 2006, Australia
 <sup>3</sup> Stellar Astrophysics Centre, Department of Physics and Astronomy, Aarhus University, Ny Munkegade 120, DK-8000 Aarhus C, Denmark
 <sup>4</sup> Department of Astronomy, The Ohio State University, Columbus, OH 43210, USA
 *Received 2018 March 12; revised 2018 April 18; accepted 2018 April 19; published 2018 May 24*

"...[dropout] is effectively a Monte Carlo integration over a Gaussian process posterior approximation (Gal & Ghahramani, 2016)."

~ Hon et al. (2018)

THE ASTROPHYSICAL JOURNAL, 859:64 (16pp), 2018 May 20 © 2018. The American Astronomical Society. All rights reserved. https://doi.org/10.3847/1538-4357/aabfdb



#### Detecting Solar-like Oscillations in Red Giants with Deep Learning

Marc Hon<sup>1</sup><sup>(1)</sup>, Dennis Stello<sup>1,2,3</sup><sup>(1)</sup>, and Joel C. Zinn<sup>4</sup><sup>(1)</sup>

<sup>1</sup> School of Physics, The University of New South Wales, Sydney NSW 2052, Australia; mtyh555@uowmail.edu.au
 <sup>2</sup> Sydney Institute for Astronomy (SIfA), School of Physics, University of Sydney, NSW 2006, Australia
 <sup>3</sup> Stellar Astrophysics Centre, Department of Physics and Astronomy, Aarhus University, Ny Munkegade 120, DK-8000 Aarhus C, Denmark
 <sup>4</sup> Department of Astronomy, The Ohio State University, Columbus, OH 43210, USA
 *Received 2018 March 12; revised 2018 April 18; accepted 2018 April 19; published 2018 May 24*



- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not restricted to variational inference).
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.
- ⊗ Principles for extension to new architectures?

- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not Dropout is a scale prior, not a posterior.
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.
- ⊗ Principles for extension to new architectures?

- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not Dropout is a scale prior, not a posterior.
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.

⊗ Principles for extension to new architectures?

 $\otimes$  Now consider the multivariate case:

$$\mathbf{h}_{l} = f_{l}(\mathbf{h}_{l-1}\boldsymbol{\Lambda}_{l}\mathbf{W}_{l}) \qquad \begin{aligned} w_{l,i,j} \sim \mathbf{N}(0,\sigma_{0}^{2}) \\ \lambda_{l,i,i} \sim p(\lambda) \end{aligned}$$

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\mathbf{\Lambda}_l\mathbf{W}_l)$$

$$w_{l,i,j} \sim \mathbf{N}(0,\sigma_0^2)$$
$$\lambda_{l,i,i} \sim p(\lambda)$$



$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\mathbf{\Lambda}_l\mathbf{W}_l)$$

$$w_{l,i,j} \sim \mathbf{N}(0,\sigma_0^2)$$
$$\lambda_{l,i,i} \sim p(\lambda)$$



$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\boldsymbol{\Lambda}_l\mathbf{W}_l)$$

$$w_{l,i,j} \sim \mathbf{N}(0,\sigma_0^2)$$
$$\lambda_{l,i,i} \sim p(\lambda)$$



⊗ Reparameterized form:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\tilde{\mathbf{W}}_l)$$



rows

#### WEIGHT MATRIX







X

⊗ Reparameterized form:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\tilde{\mathbf{W}}_l)$$

$$\tilde{w}_{l,i,j} \sim \mathbf{N}(0,\lambda_{l,i,i}^2\sigma_0^2)$$
$$\lambda_{l,i,i} \sim p(\lambda) \quad i \text{ indexes}$$

rows

Uses Bayesian shrinkage to control the effective number of hidden units.





⊗ Reparameterized form:

$$\mathbf{h}_l = f_l(\mathbf{h}_{l-1}\tilde{\mathbf{W}}_l)$$

Same structure as the **automatic relevance determination (ARD)** prior proposed by D. MacKay and R. Neal for Bayesian NNs [1994].

$$\tilde{w}_{l,i,j} \sim \mathbf{N}(0,\lambda_{l,i,i}^2\sigma_0^2)$$
$$\lambda_{l,i,i} \sim p(\lambda) \quad i \text{ indexes}$$

rows

#### WEIGHT MATRIX



## What about DropConnect?

Previous work by Wan et al. [2013] proposed *dropconnect*, which drops weights independently. Motivated by "co-adaptation" explanation of Hinton et al. [2012]



## What about DropConnect?

Previous work by Wan et al. [2013] proposed *dropconnect*, which drops weights independently. Motivated by "co-adaptation" explanation of Hinton et al. [2012]



- Revise dropout's Bayesian interpretation: should be compatible with any inference procedure (not Dropout is a scale prior, not a posterior
- Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] but has not found wide use.
- ⊗ Principles for extension to new architectures?

Revise dropout's Bayesian interpretation: should

be compatible with any inference procedure (not **Dropout is a scale** prior, not a posterior

Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] Induces ARD structure to control effective width

⊗ Principles for extension to new architectures?

⊗ Revise dropout's Bayesian interpretation: should

be compatible with any inference procedure (not <u>Dropout</u> is a scale prior, not a posterior

Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] Induces ARD structure to control effective width

⊗ Principles for extension to new architectures?













We derive a prior for ResNets called **automatic depth determination** (ADD).



Bayesian shrinkage can control the effective depth of the network

We derive a prior for ResNets called **automatic depth determination** (ADD).



Bayesian shrinkage can control the effective depth of the network

 $\tilde{w}_{l,i,i} \sim \mathsf{N}(0,\tau_l^2\sigma_0^2) \quad \tau_l \sim p(\tau)$ 

We derive a prior for ResNets called **automatic depth determination** (ADD).



Bayesian shrinkage can control the effective depth of the network

 $\tilde{w}_{l,i,j} \sim \mathbf{N}(0,\tau_l^2 \sigma_0^2) \quad \tau_l \sim p(\tau)$ 



We derive a prior for ResNets called **automatic depth determination** (ADD).



Bayesian shrinkage can control the effective depth of the network

A neural-network-analog of the *globallocal shrinkage* prior for robust regression. (See for ref: Polson & Scott [2010])



#### Extension to More Architectures



See ArXiv version (v2) of ICML paper.

Revise dropout's Bayesian interpretation: should

be compatible with any inference procedure (not **Dropout is a scale** prior, not a posterior

Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] Induces ARD structure to control effective width

⊗ Principles for extension to new architectures?

⊗ Revise dropout's Bayesian interpretation: should

be compatible with any inference procedure (not **Dropout is a scale** prior, not a posterior

Why drop hidden units?: dropping weights ('DropConnect') was explored by Wan et al. [2013] Induces ARD structure to control effective width

⊗ Principles for extension to new architectures?

Automatic depth determination for ResNets

Solution Notice The State of the State of

Solution Notice The State of the State of

- We derived two algorithms that attempt to preserve the simplicity of Dropout:
  - Tail-adaptive importance sampling for marginal-MAP inference
  - ⊗ Variational EM

Solution Notice The State of the State of

- We derived two algorithms that attempt to preserve the simplicity of Dropout:
  - Solution State State
  - ⊗ <u>Variational EM</u>
# Regression



# Regression



W



#### Test Set RMSE

ARD ADD ARD-ADD Dropout Prob. Backprop Deep GP  $2.343 \pm .31$  $2.80 \pm .13$  $2.795 \pm .16$  $2.38 \pm .12$  $\textbf{2.158} \pm .20$  $2.367 \pm .18$ Boston  $5.241 \pm .12$  $4.64 \pm .11$  $3.805 \pm .28$  $4.084 \pm .34$  $3.761 \pm .23$  $4.50 \pm .18$ Concrete  $0.852 \pm .01$  $0.867 \pm .11$  $0.853 \pm .08$  $0.47 \pm .01$  $0.903 \pm .05$  $0.57 \pm .02$ Energy  $0.064 \pm .00$  $0.071 \pm .00$  $0.066 \pm .01$  $0.064 \pm .00$ Kin8nm  $0.08 \pm .00$  $0.05 \pm .00$  $3.60 \pm .03$  $3.63 \pm .04$  $4.028 \pm .03$  $3.486 \pm .10$  $3.290 \pm .06$  $3.236 \pm .07$ Power  $0.643 \pm .01$  $0.561 \pm .03$  $0.555 \pm .01$  $0.538 \pm .03$  $0.60 \pm .01$  $0.50 \pm .01$ Wine  $0.691 \pm .12$  $0.657 \pm .14$  $0.66 \pm .06$  $0.848 \pm .05$  $0.98 \pm .09$  $0.604 \pm .16$ Yacht

# Regression



W



#### Test Set RMSE

	Dropout	Prob. Backprop	Deep GP	ARD	ADD	ARD-ADD
Boston	$2.80 \pm .13$	$2.795 \pm .16$	$2.38 \pm .12$	$2.158 \pm .20$	$2.343 \pm .31$	$2.367 \pm .18$
Concrete	$4.50 \pm .18$	$5.241 \pm .12$	$4.64 \pm .11$	$3.805 \pm .28$	$4.084 \pm .34$	$3.761 \pm .23$
Energy	$0.47 \pm .01$	$0.903 \pm .05$	$0.57 \pm .02$	$0.852 \pm .01$	$0.867 \pm .11$	$0.853 \pm .08$
Kin8nm	$0.08 \pm .00$	$0.071 \pm .00$	$0.05 \pm .00$	$0.066 \pm .01$	$0.064 \pm .00$	$0.064 \pm .00$
Power	$3.63 \pm .04$	$4.028 \pm .03$	$3.60 \pm .03$	$3.486 \pm .10$	$3.290 \pm .06$	$3.236 \pm .07$
Wine	$0.60 \pm .01$	$0.643 \pm .01$	$0.50 \pm .01$	$0.561 \pm .03$	$0.555 \pm .01$	$0.538 \pm .03$
Yacht	$0.66 \pm .06$	$0.848 \pm .05$	$0.98 \pm .09$	$0.691 \pm .12$	$0.657 \pm .14$	$0.604 \pm .16$
Avg. Rank	$4.4 \pm 1.7$	$5.6 \pm 0.5$	$3.1 \pm 1.8$	$3.0 \pm 1.1$	2.9 ±10	$2.0 \pm 1.1$

## Posterior Structure

### Heat map of summed moments (mean<sup>2</sup> + variance)



# Summary

- Clarified Dropout's modeling assumptions, generalized its Bayesian interpretation.
- Derived a new prior (ADD) to control the effective depth of Residual Networks.
- Showed our priors (w/ variational EM) can serve as direct replacement for dropout in predictive tasks.

# Thank you. Questions?



In collaboration with:



José Miguel Hernández-Lobato



Padhraic Smyth