

Hyperparameter Ensembles for Robustness and Uncertainty Quantification

Google Al

Florian Wenzelflorianwenzel@google.comJasper Snoekjsnoek@google.comDustin Trantrandustin@google.comRodolphe Jenattonrjenatton@google.com

AutoML Seminars, October 8th, 2020

Collaborators









Jasper Snoek



Dustin Tran

Context and motivation



ML systems are being deployed in many applications, e.g.,

- Conversational dialog systems, voice assistants, navigation systems
- Automatic content moderation [Link et al., 2016; Jhaver et al., 2019]

Context and motivation



ML systems are being deployed in many applications, e.g.,

- Conversational dialog systems, voice assistants, navigation systems
- Automatic content moderation [Link et al., 2016; Jhaver et al., 2019]

Those applications can be in safety-critical areas :

- Health applications [Miotto et al., 2016; Rajkomar et al., 2018; Liu et al., 2020; Mckinney et al., 2020;...]
- Self-driving cars [Levinson et al., 2011; Sun et al., 2018]
- Risk assessments [Green et al., 2019]
- Automate decisions about benefit claims & welfare issues (e.g., the Guardian, 2019)



In those critical applications, we need robust uncertainty estimation



In those critical applications, we need robust uncertainty estimation

Why?

• Knowing when to trust model's predictions, e.g., under dataset shift



In those critical applications, we need robust uncertainty estimation

Why?

- Knowing when to trust model's predictions, e.g., under dataset shift
- Better decision making, e.g., with asymmetric costs
- Active learning: Getting more data in regions where the model is uncertain
- Open set recognition
- Lifelong learning
- Exploration in reinforcement learning, Bayesian optimization, ...

Application example (1/3)

StreetView StoreFronts







Dataset shift across:

- Time
- Countries

Image source: Hendrycks et al. 2020 "The Many Faces of Robustness: A Critical Analysis of Out-of-Distribution Generalization"

Application example (2/3)

Conversational dialog systems

• Detecting out-of-scope utterances



Figure 1: Example exchanges between a user (blue, right side) and a task-driven dialog system for personal finance (grey, left side). The system correctly identifies the user's query in (1), but in (2) the user's query is mis-identified as in-scope, and the system gives an unrelated response. In (3) the user's query is correctly identified as out-of-scope and the system gives a fall-back response.

Image source: Larson et al. 2019 "An Evaluation Dataset for Intent Classification and Out-of-Scope Prediction" Google A

Application example (3/3)



• Model uncertainty to decide when to trust model vs. when to defer to human



Diabetic retinopathy detection from fundus images Gulshan et al, 2016

Eye disease classification from 3D OCT images <u>de Fauw et al, 2018</u>



But what about the behavior of modern neural networks?

Models accuracy degrades under dataset shift

Accuracy drops with

increasing shift on

Imagenet-C

Severity = 2Severity = 1Severity = 4Clean Severity = 3Severity = 50.8 0.7 0.6 0.5 Accuracy 6.0 70 8.0 8.0 Method 0.2 Dropout Vanilla Temp Scaling LL Dropout 0.1 Ensemble LL SVI 0.0 Test 2 3 Δ 5

Image source: Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift?, Ovadia et al. 2019

Google AI

Brain Team

Models accuracy degrades under dataset shift

 Accuracy drops with increasing shift on Imagenet-C



 But do the models know that they are less accurate?

Image source: Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift?, Ovadia et al. 2019

Google Al

Brain Team

Models are not calibrated under dataset shift

 Accuracy drops with increasing shift on Imagenet-C

 Calibration degrades with shift: "overconfident mistakes"



Credits: slide from B. Lakshminarayanan, D. Tran, J. Snoek "Uncertainty and Out-of-Distribution Robustness in Deep Learning"

Google AI

Brain Team



- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation



- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation

Many approaches to capture uncertainty



Various mechanisms:

- Data augmentation [Hendrycks et al., 2019; Yin et al., 2019; Thulasidasan et al., 2019]
- Architectures [Nado et al., 2020]
- Loss functions [Muller et al., 2019; Meinke et al., 2019]
- Generative models [Nalisnick et al., 2019; Ren et al., 2019]

Many approaches to capture uncertainty



Various mechanisms:

- Data augmentation [Hendrycks et al., 2019; Yin et al., 2019; Thulasidasan et al., 2019]
- Architectures [Nado et al., 2020]
- Loss functions [Muller et al., 2019; Meinke et al., 2019]
- Generative models [Nalisnick et al., 2019; Ren et al., 2019]

"Ensembles": Generate a set of diverse members

- Bayesian NNs [Hinton et al., 1993; Neal, 1995; MacKay et al., 1995; Barber et al., 1998; ...; Wenzel et al. 2020]
- Long history [Hansen et al., 1990; Levin et al., 1990; Geman et al., 1992;...; Dietterich, 2000; ...]
 - Initialization, bootstraps [Lee et al., 2015; Lakshminarayanan et al., 2017]
 - Optimization [Huang et al., 2017; Loshchilov et al., 2017; Zhang et al., 2019]
 - Hyperparameters [Caruana et al., 2004; ...; Feurer et al., 2015; Lévesque et al., 2016; Saikia et al. 2020]
 - Architectures [Zaidi+Zela et al., 2020; Antorán et al., 2020]

Many approaches to capture uncertainty



Various mechanisms:

- Data augmentation [Hendrycks et al., 2019; Yin et al., 2019; Thulasidasan et al., 2019]
- Architectures [Nado et al., 2020]
- Loss functions [Muller et al., 2019; Meinke et al., 2019]
- Generative models [Nalisnick et al., 2019; Ren et al., 2019]

"Ensembles": Generate a set of diverse members

- Bayesian NNs [Hinton et al., 1993; Neal, 1995; MacKay et al., 1995; Barber et al., 1998; ...; Wenzel et al. 2020]
- Long history [Hansen et al., 1990; Levin et al., 1990; Geman et al., 1992;...; Dietterich, 2000; ...]
 - Initialization, bootstraps [Lee et al., 2015; Lakshminarayanan et al., 2017]
 - Optimization [Huang et al., 2017; Loshchilov et al., 2017; Zhang et al., 2019]
 - Hyperparameters [Caruana et al., 2004; ...; Feurer et al., 2015; Lévesque et al., 2016; Saikia et al. 2020]
 - Architectures [Zaidi+Zela et al., 2020, Antorán et al., 2020]

Deep ensemble [Lakshminarayanan et al. 2017]



Idea:

- Multiple trainings from different seeds
- Average the predictions
- Rely on
 - Non-convexity of the objective
 - Stochasticity of initialization
 - Stochasticity of optimization
- SOTA performance [Ovadia et al., 2019; Gustafsson et al., 2019]



Inputs



Neural Network Ensembles



Deep ensembles work surprisingly well in practice





Deep ensembles are consistently among the best performing methods, especially under dataset shift

Our contribution (part 1)



- Show the importance of ensembling HPs for uncertainty modelling
- Comparison: deep ensemble vs. hyper-deep ensemble
- Simple algorithmic scheme (& little compute overhead)

Goal: Leverage diversity from both initializations and hyperparameters



- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluations

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluations



- **Random search:** Set of models with different hyperparameters
- Greedy selection: Sequentially add model minimizing ens. loss

```
Algorithm 1: greedy(\mathcal{M}, K) # Caruana et al. 2004ensemble \mathcal{E} = \{\}, score \mathcal{S}(\cdot), \mathcal{S}_{best} = +\infty;while |\mathcal{E}.unique()| \leq K dof_{\theta^*} = \arg \min_{f_{\theta} \in \mathcal{M}} \mathcal{S}(\mathcal{E} \cup \{f_{\theta}\});if \mathcal{S}(\mathcal{E} \cup \{f_{\theta^*}\}) < \mathcal{S}_{best} then| \mathcal{E} = \mathcal{E} \cup \{f_{\theta^*}\}, \mathcal{S}_{best} = \mathcal{S}(\mathcal{E});else| return \mathcal{E};endreturn \mathcal{E};
```



- **Random search:** Set of models with different hyperparameters
- Greedy selection: Sequentially add model minimizing ens. loss

```
Algorithm 1: greedy(\mathcal{M}, K) # Caruana et al. 2004ensemble \mathcal{E} = \{\}, score \mathcal{S}(\cdot), \mathcal{S}_{best} = +\infty;while |\mathcal{E}.unique()| \leq K dof_{\theta^*} = \arg \min_{f_{\theta} \in \mathcal{M}} \mathcal{S}(\mathcal{E} \cup \{f_{\theta}\});if \mathcal{S}(\mathcal{E} \cup \{f_{\theta^*}\}) < \mathcal{S}_{best} then| \mathcal{E} = \mathcal{E} \cup \{f_{\theta^*}\}, \mathcal{S}_{best} = \mathcal{S}(\mathcal{E});else| return \mathcal{E};endreturn \mathcal{E};
```



- **Random search:** Set of models with different hyperparameters
- Greedy selection: Sequentially add model minimizing ens. loss

```
Algorithm 1: greedy(\mathcal{M}, K) # Caruana et al. 2004

ensemble \mathcal{E} = \{\}, score \mathcal{S}(\cdot), \mathcal{S}_{best} = +\infty;

while |\mathcal{E}.unique()| \le K do

f_{\theta^*} = \arg \min_{f_{\theta} \in \mathcal{M}} \mathcal{S}(\mathcal{E} \cup \{f_{\theta}\});

if \mathcal{S}(\mathcal{E} \cup \{f_{\theta^*}\}) < \mathcal{S}_{best} then

| \mathcal{E} = \mathcal{E} \cup \{f_{\theta^*}\}, \mathcal{S}_{best} = \mathcal{S}(\mathcal{E});

else

| return \mathcal{E};

end

return \mathcal{E};
```



- **Random search:** Set of models with different hyperparameters
- Greedy selection: Sequentially add model minimizing ens. loss

```
Algorithm 1: greedy(\mathcal{M}, K) # Caruana et al. 2004ensemble \mathcal{E} = \{\}, score \mathcal{S}(\cdot), \mathcal{S}_{best} = +\infty;while |\mathcal{E}.unique()| \leq K dof_{\theta^*} = \arg \min_{f_{\theta} \in \mathcal{M}} \mathcal{S}(\mathcal{E} \cup \{f_{\theta}\});if \mathcal{S}(\mathcal{E} \cup \{f_{\theta^*}\}) < \mathcal{S}_{best} then| \mathcal{E} = \mathcal{E} \cup \{f_{\theta^*}\}, \mathcal{S}_{best} = \mathcal{S}(\mathcal{E});else| return \mathcal{E};endreturn \mathcal{E};
```





Hyperparameters λ





Hyperparameters $oldsymbol{\lambda}$













Concurrent work with similar construction [Zaidi+Zela et al., 2020]

Why not simply a top-K selection? [Saikia et al. 2020]





Unlike greedy, top-K is not informed by the ensemble performance



- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation

Google Al Brain Team

Settings

Small-scale experiments:

- Hyperparameters: layer-wise L2, dropout
- MLP/LeNet over CIFAR-100, Fashion MNIST

Larg(er)-scale experiments:

- Hyperparameters: block-wise L2, label smoothing
- ResNet-20 and Wide ResNet-28-10 over CIFAR-10, CIFAR-100
- Initial random search with 100 models






Results CIFAR-10







Results CIFAR-100

		single (1)	deep ens (4)	hyper-deep ens (4)
	nll ↓	$1.178\ \pm 0.020$	$\textbf{0.908}\ \pm 0.003$	$\textbf{0.896} \pm 0.003$
cifar100	acc \uparrow	0.682 ± 0.005	0.751 ± 0.002	0.754 ± 0.002
	ece ↓	0.064 ± 0.005	0.070 ± 0.005	0.050 ± 0.004
	div ↑	_	1.332 ± 0.051	1.356 ± 0.049
		single (1)	deep ens (4)	hyper-deep ens (4)
	nll↓	single (1) 0.811 ± 0.026	deep ens (4) 0.678 ± 0.013	hyper-deep ens (4) 0.636 ± 0.013
 cifar100	nll↓ acc ↑	single (1) 0.811 ± 0.026 0.801 ± 0.004	deep ens (4) 0.678 ± 0.013 0.819 ± 0.001	hyper-deep ens (4) 0.636 ± 0.013 0.831 ± 0.001
cifar100	$\begin{array}{c} \text{nll} \downarrow \\ \text{acc} \uparrow \\ \text{ece} \downarrow \end{array}$	$\begin{array}{c} \text{single (1)} \\ \hline 0.811 \ \pm \ 0.026 \\ 0.801 \ \pm \ 0.004 \\ 0.062 \ \pm \ 0.001 \end{array}$	deep ens (4) 0.678 ± 0.013 0.819 ± 0.001 0.021 ± 0.002	hyper-deep ens (4) 0.636 ± 0.013 0.831 ± 0.001 0.021 ± 0.002

≻ ResNet 20

≻ W. ResNet 28-10

Google Al Brain Team

Results CIFAR-100



Overall cost comparison



	Deep ensembles	Hyper-deep ensembles
Training cost	O(K)	O(K ²) + Cost _{random search}
Prediction cost		Same
Memory cost	S	Same

Overall cost comparison



	Deep ensembles	Hyper-deep ensembles
Training cost	O(K)	O(K ²) + Cost _{random search}
Prediction cost		Same
Memory cost		Same
"Tuning cost"	Cost' random search	0

Overall cost comparison



	Deep ensembles	Hyper-deep ensembles
Training cost	О(К)	O(K ²) + Cost _{random search}
Prediction cost		Same
Memory cost		Same
"Tuning cost"	Cost' random search	0
Total cost		\approx



Hyper-deep ensemble

- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation

For (hyper-)deep ensembles:

- Training scales linearly with number of members
- Storage scales linearly with number of members
- Prediction scales linearly with number of members





Hyperparameters λ



Google Al



Google AI

Our contribution (part 2)



- Efficient ensemble defined over different hyperparameters
- Introduce layer structure composing
 - Batch ensemble
 - Self-tuning network
- Ensemble members & their HPs are learned end-to-end in a single training

Goal:

- Exploit insights of hyper-deep ensemble...
- ...But make it more efficient!



Hyper-deep ensemble

- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation

1-slide summary about batch ensemble

Google Al Brain Tear

Batch ensemble [Wen et al., 2019]:

• Tie parameters across members: Instead of $\{\mathbf{W}_k\}_{k\in K}$

$$\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^\top) \text{ for each } k \in \{1, \dots, K\}$$



1-slide summary about batch ensemble

Batch ensemble [Wen et al., 2019]:

• Tie parameters across members: Instead of $\{\mathbf{W}_k\}_{k\in K}$

$$\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^{\top}) \text{ for each } k \in \{1, \dots, K\}$$

- Save parameters
- Fast prediction thanks to vectorization:

$$\mathbf{X}[\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^\top)] = [(\mathbf{X} \circ \mathbf{r}_k^\top) \mathbf{W}] \circ \mathbf{s}_k^\top$$

• (similar for conv. layers)





Hyper-deep ensemble

- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation

54

2-slide summary about self-tuning networks [MacKay et al. 2019]

Main idea:

$$\min_{\boldsymbol{\theta}} \left\{ \frac{1}{2} \| \mathbf{X}\boldsymbol{\theta} - \mathbf{y} \|^{2} + \lambda_{1} \| \boldsymbol{\theta} \|^{2} \right\} \quad \min_{\boldsymbol{\theta}} \left\{ \frac{1}{2} \| \mathbf{X}\boldsymbol{\theta} - \mathbf{y} \|^{2} + \lambda_{2} \| \boldsymbol{\theta} \|^{2} \right\}$$

$$\boldsymbol{\theta}^{\star}(\lambda_{1}) \qquad \boldsymbol{\theta}^{\star}(\lambda_{2})$$

• Can you track how the parameters of your model change w.r.t. HPs?



Main idea:

- Can you track how the parameters of your model change w.r.t. HPs?
- Assume you can approximate the mapping from HPs to the solution

$$\boldsymbol{\theta}^{\star}(\lambda) \approx \boldsymbol{\theta}^{\mathrm{approx}}_{\boldsymbol{\beta}}(\lambda)$$



56

2-slide summary about self-tuning networks [MacKay et al. 2019]

Main idea:

- Can you track how the parameters of your model change w.r.t. HPs?
- Assume you can approximate the mapping from HPs to the solution

$$\min_{\boldsymbol{\beta}} \mathbb{E}_{\boldsymbol{\lambda} \sim \boldsymbol{p}(\boldsymbol{\lambda})} \left\{ \frac{1}{2} \| \mathbf{X} \boldsymbol{\theta}_{\boldsymbol{\beta}}^{\text{approx}}(\boldsymbol{\lambda}) - \mathbf{y} \|^{2} + \boldsymbol{\lambda} \| \boldsymbol{\theta}_{\boldsymbol{\beta}}^{\text{approx}}(\boldsymbol{\lambda}) \|^{2} \right\} \text{ with } \boldsymbol{\theta}^{\star}(\boldsymbol{\lambda}) \approx \boldsymbol{\theta}_{\boldsymbol{\beta}}^{\text{approx}}(\boldsymbol{\lambda})$$
(to be defined)



Main idea:

- Can you track how the parameters of your model change w.r.t. HPs?
- Assume you can approximate the mapping from HPs to the solution

$$\min_{\boldsymbol{\beta}} \mathbb{E}_{\boldsymbol{\lambda} \sim \boldsymbol{p}(\boldsymbol{\lambda})} \left\{ \frac{1}{2} \| \mathbf{X} \boldsymbol{\theta}_{\boldsymbol{\beta}}^{\text{approx}}(\boldsymbol{\lambda}) - \mathbf{y} \|^{2} + \boldsymbol{\lambda} \| \boldsymbol{\theta}_{\boldsymbol{\beta}}^{\text{approx}}(\boldsymbol{\lambda}) \|^{2} \right\} \text{ with } \boldsymbol{\theta}^{\star}(\boldsymbol{\lambda}) \approx \boldsymbol{\theta}_{\boldsymbol{\beta}}^{\text{approx}}(\boldsymbol{\lambda})$$
(to be defined)

• With a single training, we could recover the entire set of solutions!



Inner workings of self-tuning networks:

• Changes captured **layer-wise**, by shifting & scaling the units: $\mathbf{W}(\boldsymbol{\lambda}) = \mathbf{W} + \Delta \circ \mathbf{e}(\boldsymbol{\lambda})^{\top}$ and $\mathbf{b}(\boldsymbol{\lambda}) = \mathbf{b} + \boldsymbol{\delta} \circ \mathbf{e}'(\boldsymbol{\lambda})$ Standard parameters



Inner workings of self-tuning networks:

• Changes captured **layer-wise**, by shifting & scaling the units: $\mathbf{W}(\boldsymbol{\lambda}) = \mathbf{W} + \Delta \circ \mathbf{e}(\boldsymbol{\lambda})^{\top} \text{ and } \mathbf{b}(\boldsymbol{\lambda}) = \mathbf{b} + \boldsymbol{\delta} \circ \mathbf{e}'(\boldsymbol{\lambda})$ Standard parameters
Embedding of HPs



Inner workings of self-tuning networks:

- Changes captured **layer-wise**, by shifting & scaling the units: $\mathbf{W}(\boldsymbol{\lambda}) = \mathbf{W} + \Delta \circ \mathbf{e}(\boldsymbol{\lambda})^{\top} \text{ and } \mathbf{b}(\boldsymbol{\lambda}) = \mathbf{b} + \boldsymbol{\delta} \circ \mathbf{e}'(\boldsymbol{\lambda})$ Standard parameters
 Embedding of HPs
- Alternating optimization:

$$\min_{\boldsymbol{\Theta}} \mathbb{E}_{\boldsymbol{\lambda} \sim p(\boldsymbol{\lambda}), (\mathbf{x}, y) \in \mathcal{D}} \big[\mathcal{L}(\mathbf{x}, y, \boldsymbol{\Theta}, \boldsymbol{\lambda}) \big]$$



Inner workings of self-tuning networks:

- Changes captured **layer-wise**, by shifting & scaling the units: $\mathbf{W}(\boldsymbol{\lambda}) = \mathbf{W} + \Delta \circ \mathbf{e}(\boldsymbol{\lambda})^{\top} \text{ and } \mathbf{b}(\boldsymbol{\lambda}) = \mathbf{b} + \boldsymbol{\delta} \circ \mathbf{e}'(\boldsymbol{\lambda})$ Standard parameters
 Embedding of HPs
- Alternating optimization:

$$\min_{\boldsymbol{\Theta}} \mathbb{E}_{\boldsymbol{\lambda} \sim p(\boldsymbol{\lambda} | \boldsymbol{\xi}_{t}), (\mathbf{x}, y) \in \mathcal{D}} \big[\mathcal{L}(\mathbf{x}, y, \boldsymbol{\Theta}, \boldsymbol{\lambda}) \big]$$

- Training step
- Update
- Keep $\boldsymbol{\xi}_t$ fixed



Inner workings of self-tuning networks:

- Changes captured **layer-wise**, by shifting & scaling the units: $\mathbf{W}(\boldsymbol{\lambda}) = \mathbf{W} + \Delta \circ \mathbf{e}(\boldsymbol{\lambda})^{\top} \text{ and } \mathbf{b}(\boldsymbol{\lambda}) = \mathbf{b} + \boldsymbol{\delta} \circ \mathbf{e}'(\boldsymbol{\lambda})$ Standard parameters
 Embedding of HPs
- Alternating optimization: $\min_{\Theta} \mathbb{E}_{\boldsymbol{\lambda} \sim p(\boldsymbol{\lambda} | \boldsymbol{\xi}_{t}), (\mathbf{x}, y) \in \mathcal{D}} \begin{bmatrix} \mathcal{L}(\mathbf{x}, y, \boldsymbol{\Theta}, \boldsymbol{\lambda}) \end{bmatrix} \qquad \min_{\boldsymbol{\xi}_{t}} \mathbb{E}_{\boldsymbol{\lambda} \sim p(\boldsymbol{\lambda} | \boldsymbol{\xi}_{t}), (\mathbf{x}, y) \in \mathcal{D}_{val}} \begin{bmatrix} \ell_{val}(f_{\boldsymbol{\Theta}}(\mathbf{x}, \boldsymbol{\lambda}), y) \end{bmatrix}$ • Training step • Update $\boldsymbol{\Theta}$ • Keep $\boldsymbol{\xi}_{t}$ fixed • Update $\boldsymbol{\xi}_{t}$



Hyper-deep ensemble

- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation

Hyper-batch ensemble (1/3)



Composite layer structure:

• Batch ensemble

 $\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^{\top})$ for each $k \in \{1, \dots, K\}$

Hyper-batch ensemble (1/3)



Composite layer structure:

• Batch ensemble

 $\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^{\top})$ for each $k \in \{1, \dots, K\}$

• Self-tuning network

$$\mathbf{W} (\boldsymbol{\lambda}) = \mathbf{W} + [\Delta] \circ \mathbf{e}(\boldsymbol{\lambda})^{\top}$$
 Model vicinity of single HP

Hyper-batch ensemble (1/3)



Composite layer structure:

• Batch ensemble

 $\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^{\top})$ for each $k \in \{1, \dots, K\}$

- Self-tuning network
 - $\mathbf{W} (\boldsymbol{\lambda}) = \mathbf{W} + [\Delta] \circ \mathbf{e}(\boldsymbol{\lambda})^{\top}$ Model vicinity of single HP
- Hyper-batch ensemble

$$\mathbf{W}_k(\boldsymbol{\lambda}_k) = \mathbf{W} \circ \overline{(\mathbf{r}_k \mathbf{s}_k^\top)} + [\Delta \circ \overline{(\mathbf{u}_k \mathbf{v}_k^\top)}] \circ \overline{\mathbf{e}(\boldsymbol{\lambda}_k)^\top} \text{ Model vicinity of } K \text{ HPs}$$

 \rightarrow Preserves memory compactness & efficient vectorization

Hyper-batch ensemble (2/3)



One HP distribution per ensemble member:

$$\min_{\boldsymbol{\Theta}} \mathbb{E}_{\boldsymbol{\Lambda}_{K} \sim q_{t}, (\mathbf{x}, y) \in \mathcal{D}} \Big[\mathcal{L}(\mathbf{x}, y, \boldsymbol{\Theta}, \boldsymbol{\Lambda}_{K}) \Big] \text{ with } \begin{bmatrix} \boldsymbol{\Lambda}_{K} = \{\boldsymbol{\lambda}_{k}\}_{k=1}^{K} \\ q_{t}(\boldsymbol{\Lambda}_{K}) = q(\boldsymbol{\Lambda}_{K} | \{\boldsymbol{\xi}_{k, t}\}_{k=1}^{K}) = \prod_{k=1}^{K} p(\boldsymbol{\lambda}_{k} | \boldsymbol{\xi}_{k, t}) \Big]$$

Hyper-batch ensemble (2/3)

One HP distribution per ensemble member:

$$\min_{\boldsymbol{\Theta}} \mathbb{E}_{\boldsymbol{\Lambda}_{K} \sim q_{t}, (\mathbf{x}, y) \in \mathcal{D}} \Big[\mathcal{L}(\mathbf{x}, y, \boldsymbol{\Theta}, \boldsymbol{\Lambda}_{K}) \Big] \text{ with } \begin{bmatrix} \boldsymbol{\Lambda}_{K} = \{\boldsymbol{\lambda}_{k}\}_{k=1}^{K} \\ q_{t}(\boldsymbol{\Lambda}_{K}) = q(\boldsymbol{\Lambda}_{K} | \{\boldsymbol{\xi}_{k, t}\}_{k=1}^{K}) = \prod_{k=1}^{K} p(\boldsymbol{\lambda}_{k} | \boldsymbol{\xi}_{k, t}) \\ p(\boldsymbol{\lambda}_{k} | \boldsymbol{\xi}_{k, t}) = q(\boldsymbol{\Lambda}_{K} | \{\boldsymbol{\xi}_{k, t}\}_{k=1}^{K}) = \prod_{k=1}^{K} p(\boldsymbol{\lambda}_{k} | \boldsymbol{\xi}_{k, t}) \Big]$$

- Which choice for $p(\boldsymbol{\lambda}_k | \boldsymbol{\xi}_{k,t})$?
 - Log-uniform parametrized by bounds
 - Popular choice for positive HPs [Bergstra et al., 2011-12]



Dynamic of λ and its lower/upper bounds 2 1 ر (*ا*المورد) المحادث λ , ensemble member 0 λ , ensemble member 1 -2 λ , ensemble member 2 -3200 400 600 800 1000 0 Epochs

Hyper-batch ensemble (3/3)



Ensemble predictions:

Individual member logits

Ensemble output

$$\begin{cases} f_{\Theta_{1}}\left(\mathbf{x},\lambda_{1}\right) \\ \vdots \\ f_{\Theta_{K}}\left(\mathbf{x},\lambda_{K}\right) \end{cases} \begin{cases} f_{\Theta}\left(\mathbf{x},\lambda\right) = \frac{1}{K}\sum_{k} \text{softmax}\left(f_{\Theta_{k}}\left(\mathbf{x},\lambda_{k}\right)\right) \end{cases}$$

 $\mathbf{W} \circ (\mathbf{r}_k \mathbf{s}_k^\top) + [\Delta \circ (\mathbf{u}_k \mathbf{v}_k^\top)] \circ \mathbf{e}(\boldsymbol{\lambda}_k)^\top$



Hyper-deep ensemble

- From deep ensemble to hyper-deep ensemble
- Construction
- Evaluation

Hyper-batch ensemble

- Two key ingredients:
 - Batch ensemble
 - Self-tuning networks
- Layer structure & end-to-end ensemble training
- Evaluation



Results CIFAR-10





Results CIFAR-100

				$\overline{}$
		batch ens (4)	hyper-batch ens (4)	
cifar100	$\begin{array}{c} nll \downarrow \\ acc \uparrow \\ ece \downarrow \\ div \uparrow \end{array}$		$\begin{array}{c} \textbf{1.152} \ \pm \ 0.015 \\ \textbf{0.699} \ \pm \ 0.002 \\ \textbf{0.095} \ \pm \ 0.002 \\ \textbf{0.159} \ \pm \ 0.007 \end{array}$	> ResNet 20
 cifar100	$\begin{array}{c} nll \downarrow \\ acc \uparrow \\ ece \downarrow \\ div \uparrow \end{array}$	batch ens (4) 0.690 ± 0.005 0.819 ± 0.001 0.026 ± 0.002 0.761 ± 0.014	hyper-batch ens (4) 0.678 ± 0.005 0.820 ± 0.000 0.022 ± 0.001 0.996 ± 0.015	W. ResNet 28-10
Results on corrupted data



W. ResNet 28-10 on CIFAR-10 Corruptions [Hendrycks et al., 2019]



Results on corrupted data



W. ResNet 28-10 on CIFAR-10 Corruptions [Hendrycks et al., 2019]



Thank you!

- Beyond continuous hyperparameters?
- Make the approach more "turnkey"
- What other types of diversity matter?
- Diversity ----- performance [Masegosa, 2019]
- Open source code:
 - Experiments in <u>uncertainty baselines</u>
 - Generic layer code in <u>edward2</u>
- Paper on Arxiv

Hyperparameter Ensembles for Robustness and Uncertainty Quantification

Florian Wenzel, Jasper Snoek, Dustin Tran, Rodolphe Jenatton Google Research {florianwenzel, jsnoek, trandustin, rjenatton}@google.com

Abstract

Ensembles over neural network weights trained from different random initialization, known as deep ensembles, achieve state-of-the-art accuracy and calibration. The recently introduced batch ensembles provide a drop-in replacement that is more parameter efficient. In this paper, we design ensembles not only over weights, but over hyperparameters to improve the state of the art in both settings. For best performance independent of budget, we propose hyper-deep ensembles, a simple procedure that involves a random search over different hyperparameters, themselves stratified across multiple random initializations. Its strong performance highlights the benefit of combining models with both weight and hyperparameter diversity. We further propose a parameter efficient version, hyper-batch ensembles, which builds on the layer structure of batch ensembles and self-tuning networks. The computational and memory costs of our method are notably lower than typical ensembles. On image classification tasks, with MLP, LeNet, and Wide ResNet 28-10 architectures, our methodology improves upon both deep and batch ensembles.

1 Introduction

Neural networks are well-suited to form *ensembles* of models [26]. Indied, neural networks trained from different random initialization can lead to equally well-performing models that are nonetheless *diverse* in that they make complementary errors on held-out data [26]. This property is explained by the multi-modal nature of their loss landscape [21] and the randomness induced by both their initialization and the stochastic methods





Additional results

EB	Google Al Brain Team
-----------	--------------------------------

		hyper-batch ens (3)	hyper-batch ens (5)	batch ens (3)	batch ens (5)	STN (1)
cifar100 (mlp)	$\begin{array}{c} \text{nll} \downarrow \\ \text{acc} \uparrow \\ \text{brier} \downarrow \\ \text{ece} \downarrow \end{array}$	$\begin{array}{c} \textbf{2.979} \pm 0.004 \\ \textbf{0.281} \pm 0.002 \\ \textbf{-0.157} \pm 0.000 \\ 0.030 \pm 0.002 \end{array}$	$\begin{array}{c} \textbf{2.983} \pm 0.001 \\ \textbf{0.282} \pm 0.001 \\ \textbf{-0.157} \pm 0.000 \\ 0.034 \pm 0.001 \end{array}$	$\begin{array}{c} 3.015 \pm 0.003 \\ 0.275 \pm 0.001 \\ \textbf{-}0.153 \pm 0.001 \\ \textbf{0.022} \pm 0.002 \end{array}$	$\begin{array}{c} 3.056 \pm 0.004 \\ 0.265 \pm 0.001 \\ \text{-}0.141 \pm 0.000 \\ 0.033 \pm 0.002 \end{array}$	$\begin{array}{c} \textbf{3.029} \pm 0.006 \\ \textbf{0.268} \pm 0.002 \\ \textbf{-0.145} \pm 0.001 \\ \textbf{0.033} \pm 0.004 \end{array}$
cifar100 (lenet)	$\begin{array}{c} \text{nll} \downarrow \\ \text{acc} \uparrow \\ \text{brier} \downarrow \\ \text{ece} \downarrow \end{array}$	$\begin{array}{c} \textbf{2.283} \pm 0.016 \\ 0.428 \pm 0.003 \\ \textbf{-0.288} \pm 0.003 \\ \textbf{0.058} \pm 0.004 \end{array}$	$\begin{array}{c} \textbf{2.297} \pm 0.009 \\ \textbf{0.425} \pm 0.002 \\ \textbf{-0.282} \pm 0.002 \\ \textbf{0.069} \pm 0.006 \end{array}$	$\begin{array}{c} \textbf{2.350} \pm 0.024 \\ \textbf{0.438} \pm 0.003 \\ \textbf{-0.295} \pm 0.003 \\ \textbf{0.058} \pm 0.015 \end{array}$	$\begin{array}{c} \textbf{2.239} \pm 0.027 \\ \textbf{0.437} \pm 0.006 \\ \textbf{-0.296} \pm 0.008 \\ \textbf{0.038} \pm 0.018 \end{array}$	$\begin{array}{c} \textbf{2.329} \pm 0.017 \\ \textbf{0.415} \pm 0.003 \\ \textbf{-0.280} \pm 0.002 \\ \textbf{0.024} \pm 0.007 \end{array}$
fmnist (mlp)	$\begin{array}{c} \text{nll} \downarrow \\ \text{acc} \uparrow \\ \text{brier} \downarrow \\ \text{ece} \downarrow \end{array}$	$\begin{array}{c} \textbf{0.308} \pm \textbf{0.002} \\ \textbf{0.892} \pm \textbf{0.001} \\ \textbf{-0.844} \pm \textbf{0.001} \\ \textbf{0.016} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} \textbf{0.304} \pm 0.001 \\ \textbf{0.892} \pm 0.001 \\ \textbf{-0.845} \pm 0.001 \\ \textbf{0.013} \pm 0.001 \end{array}$	$\begin{array}{c} 0.351 \pm 0.004 \\ 0.884 \pm 0.001 \\ \textbf{-0.830} \pm 0.001 \\ 0.020 \pm 0.001 \end{array}$	$\begin{array}{c} \textbf{0.320} \pm 0.002 \\ \textbf{0.892} \pm 0.000 \\ \textbf{-0.844} \pm 0.001 \\ \textbf{0.024} \pm 0.001 \end{array}$	$\begin{array}{c} \textbf{0.316} \pm 0.003 \\ \textbf{0.890} \pm 0.001 \\ \textbf{-0.840} \pm 0.002 \\ \textbf{0.016} \pm 0.001 \end{array}$
fmnist (lenet)	$\begin{array}{c} \text{nll} \downarrow \\ \text{acc} \uparrow \\ \text{brier} \downarrow \\ \text{ece} \downarrow \end{array}$	$\begin{array}{c} \textbf{0.212} \pm 0.001 \\ \textbf{0.924} \pm 0.001 \\ \textbf{-0.889} \pm 0.001 \\ \textbf{0.009} \pm 0.001 \end{array}$	$\begin{array}{c} \textbf{0.209} \pm 0.002 \\ \textbf{0.925} \pm 0.001 \\ \textbf{-0.891} \pm 0.001 \\ \textbf{0.008} \pm 0.001 \end{array}$	$\begin{array}{c} 0.230 \pm 0.005 \\ 0.920 \pm 0.001 \\ \text{-}0.883 \pm 0.001 \\ 0.017 \pm 0.002 \end{array}$	$\begin{array}{c} \textbf{0.221} \pm 0.002 \\ \textbf{0.922} \pm 0.001 \\ \textbf{-0.886} \pm 0.001 \\ \textbf{0.015} \pm 0.001 \end{array}$	$0.224 \pm 0.003 \\ 0.920 \pm 0.001 \\ -0.884 \pm 0.001 \\ 0.015 \pm 0.001$