

# Weight Agnostic Neural Networks

(To Appear at NeurIPS 2019)

Adam Gaier, David Ha

#### **Today's Deep Learning:**

- **Design the neural network** <u>architecture</u> (we believe to be *suitable* for encoding a policy for the task. E.g., CNNs, LSTMs, self-attention....)
- Find the <u>weight parameters</u> (for the fixed architecture, solve the *optimization* problem, using back-propagation.)

#### **Today's Deep Learning:**

- **Design the neural network** <u>architecture</u> (we believe to be *suitable* for encoding a policy for the task. E.g., CNNs, LSTMs, self-attention....)
- Find the <u>weight parameters</u> (for the fixed architecture, solve the *optimization* problem, using back-propagation.)

Question: Can we find neural architectures which are naturally capable of performing a task even when their weight parameters are randomly sampled (i.e., no learning at all)?

Question: Can we find neural architectures which are naturally capable of performing a task even when their weight parameters are randomly sampled (i.e., no learning at all)?

Question: Can we find neural architectures which are naturally capable of performing a task even when their weight parameters are randomly sampled (i.e., no learning at all)?

**Bio-plausibility (precocial behaviors):** 

- Innate ability (e.g. lizard *J* / snake 2/2 hatchlings can escape from predators..)
- Few-shot learning (e.g. ducks 🦆 are able to swim and eat, turkeys 💓 can visually recognize predators..)

Question: Can we find neural architectures which are naturally capable of performing a task even when their weight parameters are randomly sampled (i.e., no learning at all)?

**Deep learning building blocks (inductive biases):** 

- CNN, RNN, LSTM, self-attention, capsule...
- Randomly-initialized CNNs can be used effectively for image processing tasks such as superresolution, inpainting and style transfer.
- A randomly-initialized LSTM with a learned linear output layer can predict time series.

#### **Neural Architecture Search (NAS):**

- Search neural network topology using **evolutionary algorithms** (e.g., NEAT algorithm)
- Narrow the search space to architectures composed of **basic building blocks** (CNN/RNN cell/self-attention)
- Time-costly **inner loops** for training (find the optimal weight parameters)
- Never claimed that the solution is innate to the structure of the network. The **weights are the solution**; the found architectures only a substrate for the weights to inhabit.

#### **Neural Architecture Search (NAS):**

- **S** Search neural network topology using **evolutionary algorithms** (e.g., NEAT algorithm)
- D Narrow the search space to architectures composed of basic building blocks (CNN/RNN cell/self-attention)
- D Time-costly inner loops for training (find the optimal weight parameters)
- D Never claimed that the solution is innate to the structure of the network. — The weights are the solution; the found architectures only a substrate for the weights to inhabit.

#### **Bayesian Neural Networks (BNN):**

- Weight parameters are **sampled** from a distribution.
- The distribution is **learned** and usually has more parameters then the number of weights.
- Recent variance networks shows that network ensembles whose weights are sampled from a zero mean distribution can perform well on image recognition tasks (counter-intuitive...)

#### **Bayesian Neural Networks (BNN):**

- **S** Weight parameters are **sampled** from a distribution.
- D• The distribution is **learned** and usually has more parameters then the number of weights.
- **S** Recent **variance networks** shows that network ensembles whose weights are sampled from a **zero mean** distribution can perform well on image recognition tasks (counter-intuitive...)

#### **Algorithmic Information Theory (AIT):**

- Occam's razor. Simplifying the search space of its **weights**. (soft-weight sharing, regularization, etc.)
- Recent work establishes the description length of deep learning models based on **architectures**. Goal of this paper is to fined minimal architectures, instead of simple weights.

#### **Network Pruning:**

- **Pruned networks** that can achieve image classification accuracies that are much **better than chance** even with randomly initialized weights.
- **Complementary approach**: starts with a full, trained network, and takes away connections.

#### **Connectomics:**

- Wiring diagram of all neural connections of the brain. Human connectome has ~90 billion neurons and ~150 trillion synapses.
- Hope to gain insight about how the brain learns and represents memories in its **connections**.
- **Deemphasize** learning of weight parameters to test the importance of the network architecture.



**1.) Initialize** *Create population of minimal networks.*  **2.) Evaluate** Test with range of shared weight values **3.) Rank** Rank by performance and complexity

**4.) Vary** *Create new population by varying best networks.* 



**1.) Initialize** *Create population of minimal networks.* 



000

00

0







**4.) vary** Create new population by varying best networks.



- Start with a population of minimal network topologies.
- No hidden nodes. Only a fraction of the possible connections between input and output.

**1.) Initialize** *Create population oj minimal networks*. **2.) Evaluate** *Test with range of shared weight values.*  **3.) Rank** Rank by performance and complexity

**4.) Vary** *Create new population by varying best networks.* 



**1.) Initialize** *Create population of minimal networks*. **2.) Evaluate** *Test with range of shared weight values.* 



• Each network is evaluated over **multiple rollouts.** 

- Each rollout has different shared weight value assigned to all connections.
- Use fixed series of weight values (-2, -1, -0.5, +0.5, +1, +2). (assume U[-2,2] in test).
- Obtain mean performance and max performance.

**1.) Initialize** *Create population of minimal networks*. **2.) Evaluate** Test with range of shared weight values **3.) Rank** *Rank by performance and complexity* 

**4.) Vary** *Create new population by varying best networks.* 



- **Multi-objective** optimization:
  - Mean performance
  - Max performance
  - # Connections (complexity)
- No linear order: sometimes # Connections 1 but mean/max performance 1. Sorting based on dominance relationship.
- In 80% cases: ranked by (mean performance, # connections), in -Σ
  20% cases: ranked by (mean performance, max performance) -Σ

**3.) Rank** *Rank by performance and complexity* 

**4.) Vary** *Create new population by varying best networks.* 





**1.) Initialize** *Create population of minimal networks*. **2.) Evaluate** Test with range of shared weight values **3.) Rank** Rank by performance and complexity

**4.) Vary** *Create new population by varying best networks.* 







.) **Evaluate** est with range of hared weight values. **3.) Rank** Rank by performance and complexity

**4.) Vary** *Create new population by varying best networks.* 

**Possible activation functions contain:** linear, step, sin, cosine, Gaussian, tanh, sigmoid, inverse, absolute value, ReLU

negation...





#### **Three tasks:**



CartPoleSwingUp Bipedal

*BipedalWalker-v2* 

CarRacing-v0

• Inspired by biological facts: lizard 🕉, snake 🍒, ...

#### **Evaluation:**

	Wij~U(-2,2)	W~U(-2,2)	Opt.W	Opt.Wij
Swing Up	Random Weights	Random Shared Weight	Tuned Shared Weight	Tuned Weights
WANN	$57 \pm 121$	$515\pm58$	$\textbf{723} \pm \textbf{16}$	<b>932</b> ±6
Fixed Topology	$21 \pm 43$	$7\pm 2$	$8\pm1$	$918\pm7$
Biped	Random Weights	Random Shared Weight	Tuned Shared Weight	Tuned Weights
WANN	$-46\pm54$	$51\pm108$	$261 \pm 58$	$332 \pm 1$
Fixed Topology	$-129 \pm 28$	$-107 \pm 12$	$-35\pm23$	$347 \pm 1 \ [38]$
CarRacing	Random Weights	Random Shared Weight	Tuned Shared Weight	Tuned Weights
WANN	-69 $\pm$ 31	$375 \pm 177$	$\textbf{608} \pm \textbf{161}$	$893\pm74$
Fixed Topology	$-82 \pm 13$	$-85 \pm 27$	$-37 \pm 36$	<b>906 ± 21 [39]</b>

#### Performance of Randomly Sampled and Trained Weights for Continuous Control Tasks

The mean performance (over 100 trials) of the **best weight agnostic network architectures** found are compared with **standard feed forward network policies** commonly used in previous work (SOTA baselines for Biped and for CarRacing).

#### **Network Evolution:**



#### **Development of Weight Agnostic topologies over time**

*G8*: An early network which performs poorly with nearly all weights. *G32*: Relationships between the position of the cart and velocity of the pole are established. *G128*: Complexity is added to refine the balancing behavior of the elevated pole.





 $\bullet i_{\eta_V}$ 

CartpoleSwingUp champion network



BipedalWalker champion network

**Champion networks:** 



Champion network for CarRacing-v0

# **Experiment - Classification**



- The design of architectures is a focus for classification tasks. **"human-led architecture search"** on MNIST.
- High-dimensional inputs.

# **Experiment - Classification**

WANN	Test Accuracy
Random Weight	<i>82.0%</i> ± <i>18.7%</i>
Ensemble Weights	91.6%
Tuned Weight	91.9%
Trained Weights	94.2%

ANN	Test Accuracy
Linear Regression	91.6% [62]
Two-Layer CNN	99.3% [15]



#### **Classification Accuracy on MNIST**

WANNs instantiated with multiple weight values acting as an ensemble perform far better than when weights are sampled at random, and as well as a linear classifier with thousands of weights. **No single weight value has better accuracy on all digits**. The ensemble classifies samples according to the category which received the **most votes**.

#### **Experiment - Classification**



#### **Receptive Field:**

Not all neurons and connections are used to predict each digit. Starting from the output connection for a particular digit, we can map out the part of the network used to classify that digit. We can also see which parts of the inputs are used for classification.

#### Discussion



Sample shared weights from a zero-mean uniform distribution?

?

### Discussion

- Fine-tuning shared weights is useful in few-shot learning.
- **Convolutional layers** are unbeatable?
- **Baldwin effect:** not the blank slates.
- Language acquisition: Chomsky's P&P theory?

• ...