# Verification-Aided Deep Ensemble Selection

October 2022

Guy Amir Tom Zelazny Guy Katz Michael Schapira



# Deep Neural Networks (DNNs) achieve state-of-the-art results, but:



Deep Neural Networks (DNNs) achieve state-of-the-art results, but:

Training stochasticity affects accuracy and robustness



Deep Neural Networks (DNNs) achieve state-of-the-art results, but:

Training stochasticity affects accuracy and robustness

X Do not express uncertainty



Deep Neural Networks (DNNs) achieve state-of-the-art results, but:

Training stochasticity affects accuracy and robustness

X Do not express uncertainty



#### Train ensembles of DNN models



#### Train **ensembles** of DNN models

#### Higher accuracy and robustness



#### Train **ensembles** of DNN models

Higher accuracy and robustness





#### Train **ensembles** of DNN models

Higher accuracy and robustness

**Uncertainty** estimation





**Ensemble Pruning** - Given n > k trained DNNs, how can we choose the ``best'' k-sized ensemble ?

**Ensemble Pruning** - Given n > k trained DNNs, how can we choose the k-sized ensemble ?

most diverse

**Ensemble Pruning** - Given n > k trained DNNs, how can we choose the k-sized ensemble ? **most diverse** 



**Ensemble Pruning** - Given n > k trained DNNs, how can we choose the k-sized ensemble ? most diverse



#### **Option 1** – choose all k = n trained DNNs

"Ensembling Neural Networks: Many Could be Better Than All"
Zhou et. al. (2010)



use formal verification to select the best subset k



use **formal verification** to select the best subset k **How?** 

Verify all possible ensembles and choose the best one



use **formal verification** to select the best subset k **How?** *Verify all possible ensembles and choose the best one* 

**Setbacks** 



use **formal verification** to select the best subset k **How?** *Verify all possible ensembles and choose the best one* 

#### **Setbacks**

There are  $\binom{n}{k}$  k-sized ensemble combinations



use **formal verification** to select the best subset k **How?** *Verify all possible ensembles and choose the best one* 

#### **Setbacks**

There are  $\binom{n}{k}$  k-sized ensemble combinations





use **formal verification** to select the best subset k **How?** 

**Verify all possible ensembles and choose the best one** 



use **formal verification** to select the best subset k **How?** 

**Verify all possible ensembles and choose the best one** 

Check if DNNs tend to err simultaneously



Data points from the test set which are *classified the same* by all DNN members

Data points from the test set which are *classified the same* by all DNN members

Intuitively – areas with a *high consensus* among the DNNs

Data points from the test set which are *classified the same* by all DNN members

Intuitively – areas with a *high consensus* among the DNNs

Allows a *fair comparison* between DNN pairs

#### Mutual Error

Given a fixed *agreement point* and a *pair* of DNNs

A *mutual error* is a perturbation that *simultaneously* causes a pair of DNNs to *misclassify* the point

### Mutual Error (with verifier)



#### Uniqueness Score

The **mutual error score** is the average **mutual errors** a pair of DNNs have on a set of agreement points

### Uniqueness Score

The **mutual error score** is the average **mutual errors** a pair of DNNs have on a set of agreement points

The *uniqueness score* is a score indicating how (un)likely it is for a single DNN to **err** with the remaining ensemble DNNs

### Uniqueness Score

The **mutual error score** is the average **mutual errors** a pair of DNNs have on a set of agreement points

The *uniqueness score* is a score indicating how (un)likely it is for a single DNN to **err** with the remaining ensemble DNNs

The **higher** the *uniqueness score* - the **better** the DNN

**Initially:** 

1 Independently *train n DNNs* and fix a set of *agreement points* 

**Initially:** 

1 Independently *train n DNNs* and fix a set of *agreement points* 

2 Create an **arbitrary ensemble E** 

**Initially:** 

1 Independently *train n DNNs* and fix a set of *agreement points* 

2 Create an arbitrary ensemble E

#### While not TIMEOUT:

3 Compute the *uniqueness score* for each member

**Initially:** 

1 Independently *train n DNNs* and fix a set of *agreement points* 

2 Create an **arbitrary ensemble E** 

#### While not TIMEOUT:

3 Compute the *uniqueness score* for each member

4 Choose *member with lowest score* 

**Initially:** 

1 Independently *train n DNNs* and fix a set of *agreement points* 

2 Create an arbitrary ensemble E

#### While not TIMEOUT:

3 Compute the *uniqueness score* for each member

4 Choose *member with lowest score* 

5 If there is a member with a *higher* score -> *swap* 

**Initially:** 

1 Independently *train n DNNs* and fix a set of *agreement points* 

2 Create an arbitrary ensemble E

#### While not TIMEOUT:

3 Compute the *uniqueness score* for each member

4 Choose *member with lowest score* 

5 If there is a member with a *higher* score -> *swap* 









































#### But..

#### But..

#### Does it improve *robustness* on non-agreement points?



#### **Case Studies**











#### Conclusions

# A greedy-search heuristic for verification-aided ensemble selection

# Using a polynomial number of queries

Focusing on agreement points (usually) improves robustness on additional points



