

Greedily Assigning Classes to Neural Network Specialists

Sebastien Arnold

Abstract With the recent advances in deep neural networks, several experiments involved the generalist-specialist paradigm for classification. However, until now no formal study compared the performance of different clustering algorithms for class assignment. In this paper we perform such a study, suggest slight modifications to the clustering procedures, and propose a novel algorithm designed to optimize the performance of the specialist-generalist classification system. Our experiments on the CIFAR-10 and CIFAR-100 datasets allow us to investigate situations for varying number of classes on similar data. We find that our *greedy-pairs* clustering algorithm consistently outperforms other alternatives, while the choice of the confusion matrix has little impact on the final performance.

1 Introduction

Designing an efficient classification system using deep neural networks is a complicated task, which often use a multitude of models arranged in ensembles. ([3], [10]) Those ensembles often lead to state-of-the-art results on a wide range of different tasks such as image classification ([13]), speech recognition ([4]), and machine translation ([12]). The models are trained independently and in parallel, and different techniques can be used to merge their predictions.

A more structured alternative to ensembling is the use of the specialist-generalist framework. As described by [1], a natural analogy can be drawn from the medical field; a patient first consults a general practitioner who provides an initial diagnosis which is then refined by one or several specialists. In the case of classification, the doctors are replaced by neural networks and the final prediction is a combination of the specialists' outputs, and may or may not include the generalist's take.

Sebastien Arnold
Department of Computer Science, University of Southern California, e-mail: arnolds@usc.edu

In recent years, generalist and specialists have been studied under different circumstances. [5] used specialists to create an efficient image classifier for a large private dataset. The final predictions of the specialists were then used to train a reduced classifier that achieved performance similar to the whole ensemble. [6] describe a multimodal approach for emotion recognition in videos, based on specialists. Maybe closer to our work, [14] added “auxiliary heads” (acting as specialists) to their baseline network, using the precomputed features for both classification and clustering. They also underlined one of the main advantages of using specialists; a relatively low (and parallelizable) additional computational cost for increased performance.

2 Clustering Algorithms

In order to assign classes to the specialist networks, we compare several clustering algorithms on the confusion matrix of the outputs of the generalist. This confusion matrix is computed on a held-out partition of the dataset. Following previous works, we started by considering two baseline clustering algorithms, namely Lloyd’s K-Means algorithm and Spectral clustering, according to the formulation of [8]. In addition to those baseline algorithms, we evaluate the performance of two novel procedures specifically designed to improve the generalist-specialist paradigm. Those algorithms are described in the following paragraphs, and pseudo code is given in the Appendix.

We also experimented with different ways of building the confusion matrix. Besides the usual way (denoted here as *standard*) we tried three alternatives:

- *softsum*: for each prediction, we use the raw model output instead of the one-hot multi-class output,
- *softsum pred*: just like *softsum*, but only add the prediction output to the confusion matrix, if the class was correctly predicted,
- *softsum not pred*: like to *softsum pred*, but only if the prediction output was incorrectly predicted.

As discussed in later sections, the influence of the confusion matrix is minimal. Nonetheless we include them for completeness purposes.

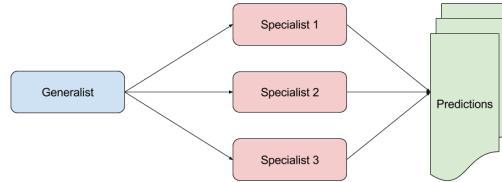


Fig. 1 An example of specialist architecture with three specialists.

Both of our clustering algorithms further modify the confusion matrix A by computing $CM = A^\top + A$, which symmetrizes the matrix. We define the entries of the matrix to be the *animosity score* between two classes; given classes a and b , their animosity score is found at $CM_{a,b}$. We then initialize each cluster with non-overlapping pairs of classes yielding maximal animosity score. Finally, we greedily select the next classes to be added to the clusters, according to the following rules:

- In the case of *greedy single* clustering, a single class maximizing the overall animosity score is added to the cluster yielding the largest averaged sum of animosity towards this class. This partitions the classes in clusters, building on the intuition that classes that are hard to distinguish should be put together.
- In the case of *greedy pairs* clustering, we follow the same strategy as in *greedy single* clustering but act on pair of classes instead of single classes. In this case we allow the clusters to overlap, and one prediction might include the opinion of several specialists.

This process is repeated until all classes have been assigned to at least one cluster.

3 Experiments

We investigate the performance of the aforementioned algorithms on the CIFAR-10 and CIFAR-100 datasets ([7]). Both datasets contain similar images, partitioned in 45'000 train, 5'000 validation, and 10'000 test images. They contain 10 and 100 classes respectively. For both experiments we train the generalist network on the train set only, and use the validation set for clustering purposes. As we are interested in the clustering performance we did not augment nor pre-process the images. Note that when trained on the horizontally flipped training and validation set our baseline algorithm reaches 10.18% and 32.22% misclassification error respectively, which is competitive with the current state-of-the-art presented in [11].

Following [2], the baseline network is based on the conclusions of [10] and uses three pairs of batch-normalized convolutional layers, each followed by a max-pooling layer, and two fully-connected layers. The same model is used for specialists, whose weights are initialized with the trained weights of the generalist.¹ One major departure from the work of [5] is that our specialists are predicting over the same classes as the generalist, i.e. we do not merge all classes outside of the cluster into a unique one. With regards to the generalist, a specialist is only biased towards a subset of the classes, since it has been fine-tuned to perform well on those ones.

¹ The code for these experiments, is freely available online at <http://www.github.com/seba-1511/specialists>.

3.1 CIFAR-10

For CIFAR-10 experiments, we considered up to five clusters, and all of the possible combinations of confusion matrix and clustering algorithms. The results for this experiments are reported in Table 1.

Table 1 Experiment results for CIFAR-10

Clustering Algorithm	standard	softsum	softsum pred	softsum not pred
spectral	(0.7046, 2)	(0.7719, 2)	(0.6989, 2)	(0.706, 2)
greedy singles	(0.5873, 2)	(0.5049, 2)	(0.5139, 3)	(0.5873, 2)
kmeans	(0.8202, 2)	(0.8202, 2)	(0.8202, 2)	(0.8202, 2)
greedy pairs	(0.8835, 2)	(0.8835, 2)	(0.8727, 3)	(0.8835, 2)

Interestingly, the choice of confusion matrix has only a limited impact on the overall performance, indicating that the emphasis should be put on the clustering algorithm. We notice that clustering with greedy pairs consistently yields better scores. However none of the specialist experiments is able to improve on the baseline, suggesting that specialists might not be the framework of choice when dealing with a small number of classes.

3.2 CIFAR-100

For CIFAR-100 we performed the exact same experiment as for CIFAR-10 but used more specialists, the largest experiments involving 28 clusters. The results are shown in Table 2.

Table 2 Experiment results for CIFAR-100

Clustering Algorithm	standard	softsum	softsum pred	softsum not pred
spectral	(0.5828, 2)	(0.5713, 2)	(0.5755, 2)	(0.5795, 3)
greedy singles	(0.3834, 2)	(0.3733, 2)	(0.3803, 2)	(0.3551, 2)
kmeans	(0.5908, 2)	(0.5618, 2)	(0.5820, 3)	(0.5876, 2)
greedy pairs	(0.6141, 6)	(0.5993, 6)	(0.6111, 6)	(0.607, 6)

Similarly to CIFAR-10, we observe that greedy pairs clustering outperforms the other clustering techniques, and that the different types of confusion matrix have a limited influence on the final score. We also notice that fewer clusters tend to work better. Finally, and unlike the results for CIFAR-10, some of the specialists are able to improve upon the generalist, which confirms our intuition that specialists are better suited to problems involving numerous output classes.

We suggest the following explanation for the improved performance of greedy pairs is the following. Allowing clusters to overlap leads to the assignment of difficult classes to multiple specialists. At inference time, more networks will influence the final prediction which is analogous to building a larger ensemble for difficult classes.

4 Conclusion and Future Work

We introduced a novel clustering algorithm for the specialist-generalist framework, which is able to consistently outperform other techniques. We also provided a preliminary study of the different factors coming into play when dealing with specialists, and concluded that the choice of confusion matrix from our proposed set only has little impact on the final classification outcome.

Despite our encouraging results with clustering techniques, no one of our specialists-based experiments came close to compete with the generalist model trained on the entire train and validation set. This was a surprising outcome and we suppose that this effect comes from the size of the datasets. In both cases, 5'000 images corresponds to 10% of the original training set and removing that many training examples has a drastic effect on both generalists and specialists. All the more so since we are not using any kind of data augmentation techniques, which could have moderated this downside. An obvious future step is to validate the presented ideas on a much larger dataset such as [9] where splitting the train set would not hurt the train score as much.

Acknowledgements We would like to thank Greg Ver Steeg, Gabriel Pereyra, and Pranav Rajpurkar for their comments and advices. We also thank Nervana Systems for providing GPUs as well as their help with their deep learning framework.

5 Appendix

Algorithm 1 Greedy Pairs Clustering

```

1: procedure GREEDYPAIRS( $M, N$ ) ▷ Confusion matrix  $M$ , number of clusters  $N$ 
2:    $M \leftarrow M + M^T$ 
3:   Initialize  $N$  clusters with non-overlapping pairs maximizing the entries of  $M$ .
4:   while every class has not been assigned do
5:     Get the next pair  $(a, b)$  maximizing the entry in  $M$ 
6:     cluster =  $\text{cclustersargmin}(\text{Animosity}(a, c) + \text{Animosity}(b, c))$ 
7:     Assign(cluster,  $a, b$ )
8:   return clusters

```

References

1. Bochereau, Laurent, and Bourguine, Paul. A Generalist-Specialist Paradigm for Multilayer Neural Networks. *Neural Networks*, 1990.
2. Courbariaux, Matthieu, Bengio, Yoshua, and David, Jean-Pierre. BinaryConnect: Training Deep Neural Networks with Binary Weights during Propagations. *NIPS*, 2015.
3. Dieleman, Sander, Willett, Kyle W., and Dambre, Joni. Rotation-invariant convolutional neural networks for galaxy morphology prediction. *Oxford Journals*, 2015.
4. Hannun, Awni, Case, Carl, Casper, Jared, Catanzaro, Bryan, Diamos, Greg, Elsen, Erich, Prenger, Ryan, Satheesh, Sanjeev, Sengupta, Shubho, Coates, Adam, and Ng, Andrew Y. Deep Speech: Scaling up end-to-end speech recognition. *Arxiv Preprint*, 2014.
5. Hinton, Geoffrey E., Vinyals, Oriol, and Dean, Jeff. Distilling the Knowledge in a Neural Network. *NIPS 2014 Deep Learning Workshop*.
6. Kahou, Samira Ebrahimi, Bouthiller, Xavier, Lamblin, Pascal, Gulcehre, Caglar, Michalski, Vincent, Konda, Kishore, Jean, Sbastien, Froumenty, Pierre, Dauphin, Yann, Boulanger-Lewandowski, Nicolas, Ferrari, Raul Chandias, Mirza, Mehdi, Warde-Farley, David, Courville, Aaron, Vincent, Pascal, Memisevic, Roland, Pal, Christopher, and Bengio, Yoshua. EmoNets: Multimodal deep learning approaches for emotion recognition in video. *Journal on Multimodal User Interfaces*, 2015.
7. Krizhevsky, Alex. Learning Multiple Layers of Features from Tiny Images. 2009.
8. Ng, Andrew Y., Jordan, Michael I., Weiss, Yair. On spectral clustering: Analysis and an algorithm. *NIPS* 2002.
9. Russakovsky, Olga, Deng, Jia, Su, Hao, Krause, Jonathan, Satheesh, Sanjeev, Ma, Sean, Huang, Zhiheng, Karpathy, Andrej, Khosla, Aditya, Bernstein, Michael, Berg, Alexander C., and Fei-Fei, Li. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 2015.
10. Simonyan, Karen and Zisserman, Andrew. Very Deep Convolutional Networks for Large-Scale Image Recognition. *International Conference on Learning Representations*, 2015.
11. Springenberg, Jost Tobias, Dosovitskiy, Alexey, Brox, Thomas, and Riedmiller, Martin. Striving for Simplicity: The All Convolutional Net. *International Conference on Learning Representations Workshop*, 2015.
12. Sutskever, Ilya, Vinyals, Oriol, and Le, Quoc V. Sequence to Sequence Learning with Neural Networks. *Arxiv Preprint*, 2014.
13. Szegedy, Christian, Liu, Wei, Jia, Yangqing, Sermanet, Pierre, Reed, Scott, Anguelov, Dragomir, Erhan, Dumitru, Vanhoucke, Vincent, and Rabinovich, Andrew. Going deeper with convolutions. *Arxiv Preprint*, 2014.
14. Warde-Farley, David, Rabinovich, Andrew, and Anguelov, Dragomir. Self-Informed Neural Networks Structure Learning. *International Conference on Representations Learning*, 2015.