

# Deep Neural Decision Forests

Peter Kotschieder<sup>1</sup>   Madalina Fiterau<sup>2</sup>   Antonio Criminisi<sup>1</sup>   Samuel Rota Bulò<sup>1,3</sup>

Microsoft Research<sup>1</sup>   Carnegie Mellon University<sup>2</sup>   Fondazione Bruno Kessler<sup>3</sup>  
Cambridge, UK   Pittsburgh, PA   Trento, Italy

## Supplementary Material

This document provides the following supplementary contributions:

- Additional details regarding the ImageNet experiment (see Section 1).
- Detailed derivations for the gradient term depending on the decision tree splits (see Section 2).
- The proof that our update rule for the leaf predictions  $\pi$  in Equ. (11) of our main ICCV paper monotonically decreases the risk  $R$  until a fixed point is reached (see Section 3).

### 1. ImageNet experiment: GoogLeNet vs. dNDF.NET architectures

This section provides additional description for the ImageNet experiment [3]. In particular, we describe and illustrate the changes we made to GoogLeNet [4] to obtain our proposed dNDF.NET architecture, using deep neural decision forests (dNDFs) as classifiers. The GoogLeNet architecture we have used as basis for our experiments (see left illustration in Fig. 1, taken from [4]) has a reported Top5-Error of 10.07%, when used in a single-model, single-crop setting (see first row in Tab. 3 of [4]).

Fig. 1 (right illustration) shows that we have introduced two different modifications with respect to the original GoogLeNet architecture. First, we have connected the outputs of the Concat layers to the inputs of the AveragePool layers (as described in the main paper), visualized by red arrows in the plot. The resulting, modified baseline network is dubbed GoogLeNet $\star$  and achieves a Top5-Error of 10.02% when using conventional SoftMax layers as in the original network. The implementation yielding this score was realized in the Distributed (Deep) Machine Learning Common (DMLC) library [2, 1]<sup>1</sup>, using resized images with dimensionality 100x100 as described in [2]. The training used the standard settings for GoogLeNet, stochastic gradient descent with 0.9 momentum, fixed learning rate schedule, decreasing the learning rate by 4% every 8 epochs. We trained GoogLeNet $\star$  with mini-batches composed of 50 images.

In order to obtain a *Deep Neural Decision Forest* architecture coined dNDF.NET, we have replaced each Softmax layer from GoogLeNet $\star$  with a forest consisting of 10 trees (each fixed to depth 15), resulting in a total number of 30 trees. For our architecture, which we implemented in DMLC as well, we trained the network for 1000 epochs using mini-batches composed of 100.000 images. This is feasible due to distribution of the computational load to a cluster of 52 CPUs and 12 hosts, where each host is equipped with a NVIDIA Tesla K40 GPU.

We refer to the individual forests as dNDF<sub>0</sub>, dNDF<sub>1</sub> and dNDF<sub>2</sub>, where dNDF<sub>0</sub> is closest to the input layer and dNDF<sub>2</sub> is the final (last) layer in the architecture. Each tree is a balanced and fixed depth 15 tree, which means that the total number of per-tree split nodes is  $2^{15} - 1 = 32.767$  and the number of leaf nodes is  $2^{15} = 32.768$ .

---

<sup>1</sup><https://github.com/dmlc/cxxnet.git>

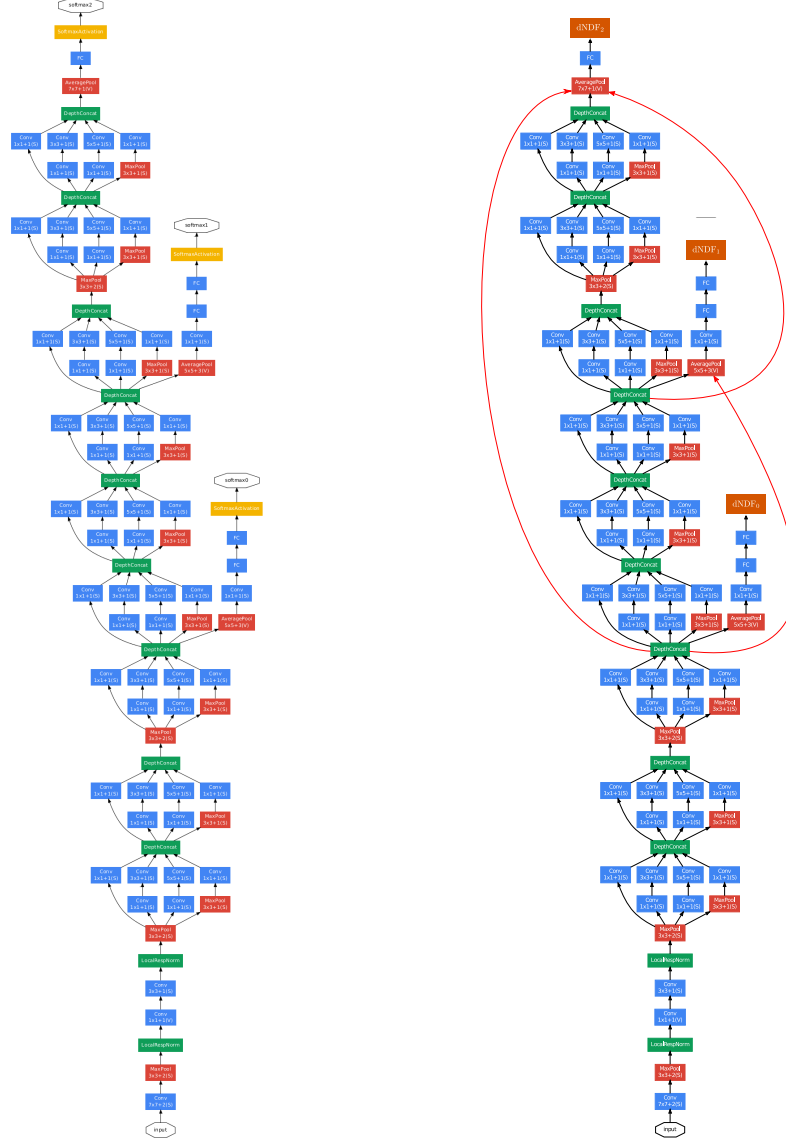


Figure 1. Left: Original GoogLeNet architecture proposed in [4]. Right: The modifications we brought to the GoogLeNet architecture resulting in dNDF.NET – our proposed model using dNDFs as final classifiers. Best viewed with digital zoom.

## 2. Split function gradient term derivations

Given the definitions for our split decision functions  $d_n(\mathbf{x}; \Theta)$  and the log-loss  $L(\Theta, \pi; \mathbf{x}, y)$  (see Equ. (3) and Equ. (6) in the main paper, respectively), we can derive the gradient term in Equ. (9) of the main paper as follows:

$$\begin{aligned}
 \frac{\partial L(\Theta, \pi; \mathbf{x}, y)}{\partial f_n(\mathbf{x}; \Theta)} &= \sum_{\ell \in \mathcal{L}} \frac{\partial L(\Theta, \pi; \mathbf{x}, y)}{\partial \mu_\ell(\mathbf{x}; \Theta)} \frac{\partial \mu_\ell(\mathbf{x}; \Theta)}{\partial f_n(\mathbf{x}; \Theta)} \\
 &= - \sum_{\ell \in \mathcal{L}} \frac{\pi_{\ell y}}{\mathbb{P}_T[y|\mathbf{x}, \Theta, \pi]} \frac{\partial \mu_\ell(\mathbf{x}; \Theta)}{\partial f_n(\mathbf{x}; \Theta)} \\
 &= - \sum_{\ell \in \mathcal{L}} \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}, \Theta, \pi]} \frac{\partial \log \mu_\ell(\mathbf{x}; \Theta)}{\partial f_n(\mathbf{x}; \Theta)},
 \end{aligned}$$

where

$$\begin{aligned}\frac{\partial \log \mu_\ell(\mathbf{x}; \Theta)}{\partial f_n(\mathbf{x}; \Theta)} &= \mathbb{1}_{\ell \prec n} \frac{\partial \log d_n(\mathbf{x}; \Theta)}{\partial f_n(\mathbf{x}; \Theta)} \\ &\quad + \mathbb{1}_{n \searrow \ell} \frac{\partial \log \bar{d}_n(\mathbf{x}; \Theta)}{\partial f_n(\mathbf{x}; \Theta)} \\ &= \mathbb{1}_{\ell \prec n} \bar{d}_n(\mathbf{x}; \Theta) - \mathbb{1}_{n \searrow \ell} d_n(\mathbf{x}; \Theta).\end{aligned}$$

By substituting the latter in the previous formula we get

$$\begin{aligned}\frac{\partial L(\Theta, \boldsymbol{\pi}; \mathbf{x}, y)}{\partial f_n(\mathbf{x}; \Theta)} &= - \sum_{\ell \in \mathcal{L}} \mathbb{1}_{\ell \prec n} \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}, \Theta, \boldsymbol{\pi}]} \bar{d}_n(\mathbf{x}; \Theta) \\ &\quad + \sum_{\ell \in \mathcal{L}} \mathbb{1}_{n \searrow \ell} \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}, \Theta, \boldsymbol{\pi}]} d_n(\mathbf{x}; \Theta) \\ &= - \sum_{\ell \in \mathcal{L}_{n_l}} \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}, \Theta, \boldsymbol{\pi}]} \bar{d}_n(\mathbf{x}; \Theta) \\ &\quad + \sum_{\ell \in \mathcal{L}_{n_r}} \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}, \Theta, \boldsymbol{\pi}]} d_n(\mathbf{x}; \Theta) \\ &= d_n(\mathbf{x}; \Theta) A_{n_r} - \bar{d}_n(\mathbf{x}; \Theta) A_{n_l}.\end{aligned}$$

### 3. Proof of update rule for $\boldsymbol{\pi}$

**Theorem 1.** Consider a tree with parameters  $\Theta$  and  $\boldsymbol{\pi}$  and let

$$\hat{\pi}_{\ell y} = \frac{1}{Z_\ell} \sum_{(\mathbf{x}, y') \in \mathcal{T}} \mathbb{1}_{y=y'} \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}; \Theta, \boldsymbol{\pi}]}, \quad \text{for all } (\ell, y) \in \mathcal{L} \times \mathcal{Y}, \quad (12)$$

where  $Z_\ell$  is the normalizing factor ensuring that  $\hat{\boldsymbol{\pi}}_\ell = (\hat{\pi}_{\ell y})_{y \in \mathcal{Y}}$  is a probability distribution. In other terms, we assume  $\hat{\pi}_{\ell y}$  to be the result of an update step as per (11) of our ICCV contribution. The following holds:

$$R(\Theta, \boldsymbol{\pi}; \mathcal{T}) \geq R(\Theta, \hat{\boldsymbol{\pi}}; \mathcal{T})$$

with equality if and only if  $\hat{\boldsymbol{\pi}} = \boldsymbol{\pi}$ , where  $R$  is the risk defined in (5) of our ICCV contribution, and  $\boldsymbol{\pi} = (\pi_\ell)_{\ell \in \mathcal{L}}$ .

*Proof.* Consider the following auxiliary function:

$$\phi(\boldsymbol{\pi}, \bar{\boldsymbol{\pi}}) = R(\Theta, \bar{\boldsymbol{\pi}}; \mathcal{T}) - \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \sum_{\ell \in \mathcal{L}} \xi_\ell(\bar{\boldsymbol{\pi}}; \mathbf{x}, y) \log \left( \frac{\pi_{\ell y}}{\bar{\pi}_{\ell y}} \right),$$

where

$$\xi_\ell(\boldsymbol{\pi}; \mathbf{x}, y) = \frac{\pi_{\ell y} \mu_\ell(\mathbf{x}; \Theta)}{\mathbb{P}_T[y|\mathbf{x}; \Theta, \boldsymbol{\pi}]}.$$

and  $\mathbb{P}_T[y|\mathbf{x}, \Theta, \boldsymbol{\pi}]$  is defined as per (1) of our ICCV contribution. Note that  $\phi(\boldsymbol{\pi}, \boldsymbol{\pi}) = R(\Theta, \boldsymbol{\pi}; \mathcal{T})$  holds for any  $\boldsymbol{\pi}$ , for the logarithm term in  $\phi$  nullifies. Moreover,  $\phi(\boldsymbol{\pi}, \bar{\boldsymbol{\pi}}) \geq R(\Theta, \boldsymbol{\pi}; \mathcal{T})$  holds for any  $\boldsymbol{\pi}$  and  $\bar{\boldsymbol{\pi}}$ . This can be seen by applying

Jensen's inequality and with few algebraic manipulations:

$$\begin{aligned}
R(\Theta, \pi; \mathcal{T}) &= -\frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \log \left( \sum_{\ell \in \mathcal{L}} \pi_{\ell y} \mu_{\ell}(\mathbf{x}; \Theta) \right) \\
&\leq -\frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \sum_{\ell \in \mathcal{L}} \xi_{\ell}(\bar{\pi}; \mathbf{x}, y) \log \left( \frac{\pi_{\ell y} \mu_{\ell}(\mathbf{x}; \Theta)}{\xi_{\ell}(\bar{\pi}; \mathbf{x}, y)} \right) \quad (\text{by Jensen's inequality}) \\
&= -\frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \sum_{\ell \in \mathcal{L}} \xi_{\ell}(\bar{\pi}; \mathbf{x}, y) \left[ \log \left( \frac{\pi_{\ell y}}{\bar{\pi}_{\ell y}} \right) + \log \mathbb{P}_T[y|\mathbf{x}, \Theta, \bar{\pi}] \right] \\
&= R(\Theta, \bar{\pi}; \mathcal{T}) - \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \sum_{\ell \in \mathcal{L}} \xi_{\ell}(\bar{\pi}; \mathbf{x}, y) \log \left( \frac{\pi_{\ell y}}{\bar{\pi}_{\ell y}} \right) = \phi(\pi, \bar{\pi}).
\end{aligned}$$

We can now show that  $\hat{\pi}$  is a global minimizer of  $\phi(\cdot, \pi)$  for any value of  $\pi$ . We start rewriting  $\hat{\pi}$  in terms of  $\xi_{\ell}$  as follows:

$$\hat{\pi}_{\ell y} = \frac{1}{Z_{\ell}} \sum_{(\mathbf{x}, y') \in \mathcal{T}} \mathbb{1}_{y=y'} \xi_{\ell}(\pi; \mathbf{x}, y'),$$

where  $Z_{\ell}$  is the normalizing factor. Then, we have that

$$\begin{aligned}
\phi(\hat{\pi}, \pi) - \phi(\bar{\pi}, \pi) &= -\frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \sum_{\ell \in \mathcal{L}} \xi_{\ell}(\pi; \mathbf{x}, y) \log \left( \frac{\hat{\pi}_{\ell y}}{\bar{\pi}_{\ell y}} \right) \\
&= -\frac{1}{|\mathcal{T}|} \sum_{\ell \in \mathcal{L}} \sum_{y \in \mathcal{Y}} \sum_{(\mathbf{x}, y') \in \mathcal{T}} \mathbb{1}_{y=y'} \xi_{\ell}(\pi; \mathbf{x}, y) \log \left( \frac{\hat{\pi}_{\ell y}}{\bar{\pi}_{\ell y}} \right) \\
&= -\frac{1}{|\mathcal{T}|} \sum_{\ell \in \mathcal{L}} \sum_{y \in \mathcal{Y}} Z_{\ell} \hat{\pi}_{\ell y} \log \left( \frac{\hat{\pi}_{\ell y}}{\bar{\pi}_{\ell y}} \right) \\
&= -\frac{1}{|\mathcal{T}|} \sum_{\ell \in \mathcal{L}} Z_{\ell} D_{KL}(\hat{\pi}_{\ell} \| \bar{\pi}_{\ell}) \leq 0,
\end{aligned}$$

holds for all values of  $\bar{\pi}$ , where  $D_{KL}(\cdot \| \cdot)$  is the Kullback-Leibler divergence. Note that the last inequality yields equality if and only if  $\hat{\pi} = \bar{\pi}$ , for the Kullback-Leibler divergence yields zero if and only if the two distributions in input coincide. Accordingly,  $\hat{\pi}$  is a *strict* global minimizer of  $\phi(\cdot, \pi)$  for any  $\pi$ .

As a consequence of the previous derivations we have

$$R(\Theta, \pi; \mathcal{T}) = \phi(\pi, \pi) > \phi(\hat{\pi}, \pi) \geq R(\Theta, \hat{\pi}; \mathcal{T}),$$

where equality holds if and only if we have a fixed point of the update rule (12), *i.e.* if  $\hat{\pi} = \pi$ . □

## References

- [1] M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B.-Y. Su. Scaling distributed machine learning with the parameter server. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. USENIX Association, 2014. [1](#)
- [2] M. Li, L. Zhou, Z. Yang, A. Li, F. Xia, D. G. Andersen, and A. Smola. Parameter server for distributed machine learning. In *Big Learning NIPS Workshop*, 2013. [1](#)
- [3] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 2014. [1](#)
- [4] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014. [1, 2](#)