

# Towards Hierarchical Explanation

Christiaan   Hinrik   Albert   Anna

FACT-AI 2020

# Table of Contents

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- 1 Reproducing the prototype network
- 2 The hierarchical prototype network
- 3 Hierarchical results
- 4 Discussion & Broader implications

# Original paper

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications



Li, Oscar and Liu, Hao and Chen, Chaofan and Rudin, Cynthia.

Deep learning for case-based reasoning through prototypes:  
A neural network that explains its predictions.

*Thirty-Second AAAI Conference on Artificial Intelligence,*  
2018

# Idea

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- Broadly speaking, (convolutional) neural nets are not interpretable
- Instead of explaining predictions after training, integrate explanations in training goal
- Learn a fixed amount of **prototypes** which *represent the entire dataset*

# The prototype network

Towards  
Hierarchical  
Explanation

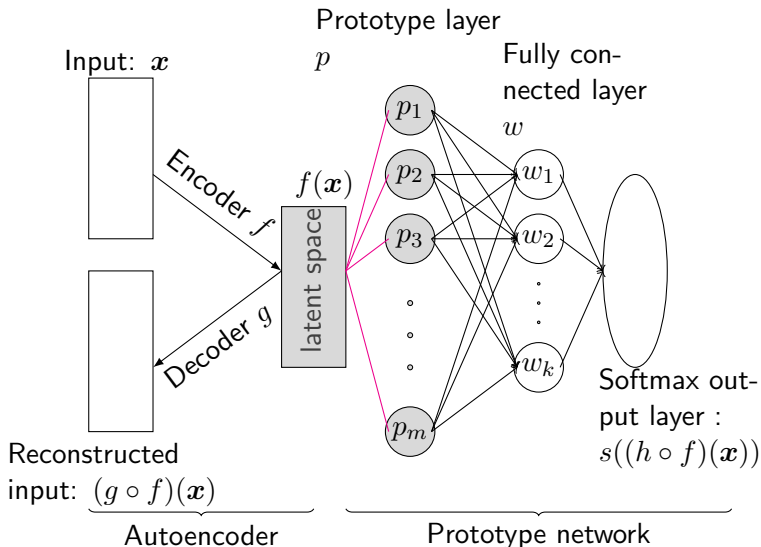
Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications



# Building the loss function

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

## Loss function

Loss = Reconstruction error

$$L((f, g), D) = R(g \circ f, D)$$

# Building the loss function

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

## Loss function

Loss = Crossentropy loss + Reconstruction error

$$L((f, g, h), D) = E(h \circ f, D) + R(g \circ f, D)$$

# Building the loss function

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

## Loss function

Loss = Crossentropy loss + Reconstruction error +  
Regularization terms

$$L((f, g, h), D) = E(h \circ f, D) + R(g \circ f, D) + R_1 + R_2$$



# Building the loss function

## Loss function

Loss = Crossentropy loss + Reconstruction error +  
Regularization terms

$$L((f, g, h), D) = E(h \circ f, D) + R(g \circ f, D) + R_1 + R_2$$

## Regularization terms for prototypes $p_1, \dots, p_m$

$$R_1(p_1, p_2, \dots, p_m, D) = \frac{1}{m} \sum_{j=1}^m \min_{i \in [1, n]} \|p_j - f(x_i)\|_2^2$$

$$R_2(p_1, p_2, \dots, p_m, D) = \frac{1}{n} \sum_{i=1}^n \min_{j \in [1, m]} \|p_j - f(x_i)\|_2^2$$

# Building the loss function

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

## Loss function

Loss = Crossentropy loss + Reconstruction error +  
Regularization terms

$$L((f, g, h), D) = \lambda_{\text{class}} E(h \circ f, D) + \lambda_R R(g \circ f, D) + \lambda_1 R_1 + \lambda_2 R_2$$

## Regularization terms for prototypes $p_1, \dots, p_m$

$$R_1(p_1, p_2, \dots, p_m, D) = \frac{1}{m} \sum_{j=1}^m \min_{i \in [1, n]} \|p_j - f(x_i)\|_2^2$$

$$R_2(p_1, p_2, \dots, p_m, D) = \frac{1}{n} \sum_{i=1}^n \min_{j \in [1, m]} \|p_j - f(x_i)\|_2^2$$

# Reproducing results

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- MNIST digits
- 15 prototypes

# Reproducing results

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- MNIST digits
- 15 prototypes
- Autoencoder with four convolutional layers

# Reproducing results

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- MNIST digits
- 15 prototypes
- Autoencoder with four convolutional layers
- Learning rate 0.0001, Epochs 1500

# Reproducing results

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- MNIST digits
- 15 prototypes
- Autoencoder with four convolutional layers
- Learning rate 0.0001, Epochs 1500
- Test accuracy: 98.879% (Paper reports 99.22%)

# Learned prototypes

- Original results:



Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

# Learned prototypes

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

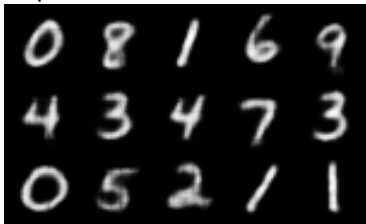
Hierarchical  
results

Discussion &  
Broader  
implications

- Original results:



- Reproduced results:





# However...

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

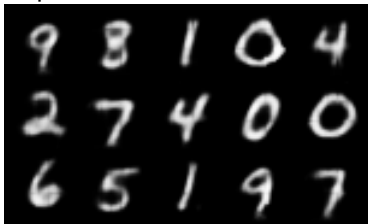
Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- Reproduced results with another seed:



- (Accuracy still 98.71%)

# The hierarchical idea

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- If  $m > K$ , multiple prototypes of 1 class
- Sometimes prototype network does not learn a prototype for each class
- If  $m = K$ , 1 prototype for each class
- Cannot capture intraclass differences

# The hierarchical idea

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- If  $m > K$ , multiple prototypes of 1 class
- Sometimes prototype network does not learn a prototype for each class
- If  $m = K$ , 1 prototype for each class
- Cannot capture intraclass differences

Possible solution: *superprototypes*

Input example  $\prec$  Subprototype  $\prec$  Superprototype,  
where  $\prec$  means “more specific than”

# Our architecture

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

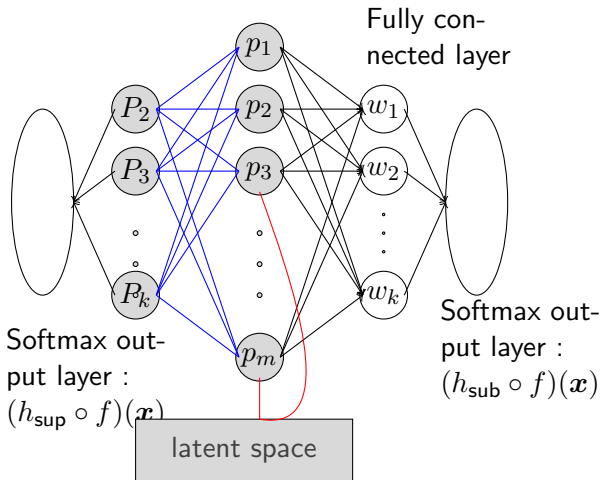
Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

Superprototype layer  $P$     Subprototype network  $p$



# Two new loss terms

## New loss term

$$L((f, g, h), D) = \lambda_{\text{class}} E(h \circ f, D) + \lambda_R R(g \circ f, D) + \lambda_1 R_1 + \lambda_2 R_2 + \lambda_3 R_3 + \lambda_4 R_4$$

$$R_1(\mathbf{p}_1, \dots, \mathbf{p}_m, D) = \frac{1}{m} \sum_{j=1}^m \min_{i \in [1, n]} \|\mathbf{p}_j - f(\mathbf{x}_i)\|_2^2$$

$$R_2(\mathbf{p}_1, \dots, \mathbf{p}_m, D) = \frac{1}{n} \sum_{i=1}^n \min_{j \in [1, m]} \|\mathbf{p}_j - f(\mathbf{x}_i)\|_2^2$$

$$R_3(\mathbf{P}_1, \dots, \mathbf{P}_K, \mathbf{p}_1, \dots, \mathbf{p}_m) = \frac{1}{K} \sum_{k=1}^K \min_{j \in [1, m]} \|\mathbf{P}_k - \mathbf{p}_j\|_2^2$$

$$R_4(\mathbf{P}_1, \dots, \mathbf{P}_K, \mathbf{p}_1, \dots, \mathbf{p}_m) = \frac{1}{m} \sum_{j=1}^m \min_{k \in [1, K]} \|\mathbf{P}_k - \mathbf{p}_j\|_2^2$$

# Results

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- Accuracy for superprototype classifier: 98.86%
- Accuracy for subprototype classifier: 99.02%

# Superprototypes and subprototypes

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

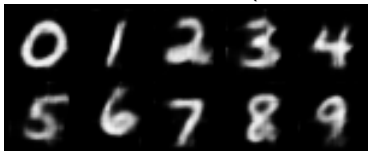
Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

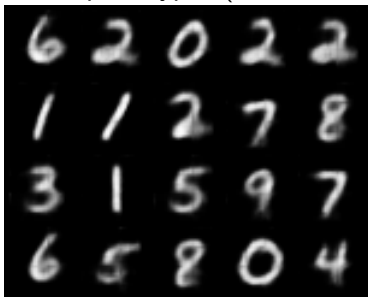
Hierarchical  
results

Discussion &  
Broader  
implications

- $K$  Superprototypes (fixed layer)



- $m$  Subprototypes (learnable FC layer)



# Transparency & Fairness

Towards  
Hierarchical  
Explanation

Christiaan,  
Hinrik, Albert,  
Anna

Reproducing  
the prototype  
network

The  
hierarchical  
prototype  
network

Hierarchical  
results

Discussion &  
Broader  
implications

- Model interclass and intraclass variation
- Some hierarchical interpretability
- Possibly discover biases in dataset by looking at (sub)prototypes



Thank you for your attention

Thank you for your attention  
Any questions?