
Deep Learning for Case-Based Reasoning through Prototypes

Oscar Li , Hao Liu , Chaofan Chen , Cynthia Rudin

(Reimplementation)

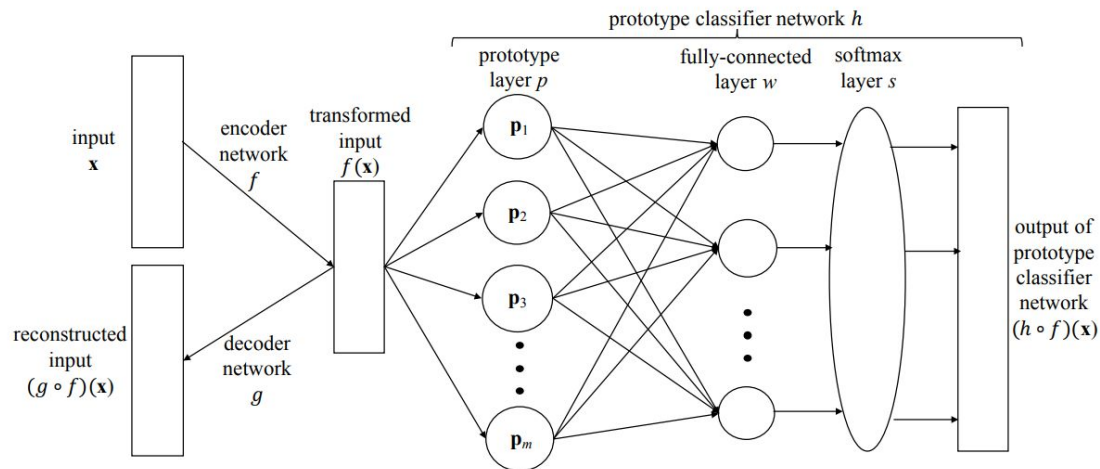
Main idea

Add special prototype layer to learn prototypes.

Use prototypes to explain class prediction for test examples.

	8	9	0	7	3
	0.98	1.47	0.70	1.55	1.49
6	0.29	1.69	1.02	0.41	0.15
	5	2	2	4	2
	0.88	1.40	1.45	1.28	1.28

Prototypes similarity



Desired prototypes are:

$$R(g \circ f, D) = \frac{1}{n} \sum_{i=1}^n \|(g \circ f)(\mathbf{x}_i) - \mathbf{x}_i\|_2^2.$$

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

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

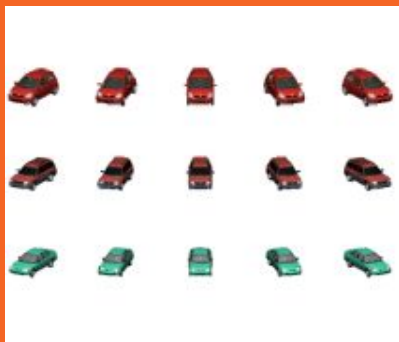
- well reconstructed (minimizing autoencoder error)
- as close as possible to at least one of the training examples in the latent space
- representative for all training examples (i.e. every encoded training example is as close as possible to one of the prototype vectors)

*We also use cross-entropy to penalize misclassification



MNIST

CARS



Datasets



FASHION-MNIST

ROCK, PAPER, SCISSORS



Learned prototypes



Distances between prototypes and test example

	0.98	1.47	0.70	1.55	1.49
	0.29	1.69	1.02	0.41	0.15
	0.88	1.40	1.45	1.28	1.28



	2.2	2.87	1.85	0.1	3.53
	2.67	2.39	3.07	1.91	1.59
	1.54	2.21	2.52	1.65	2.79

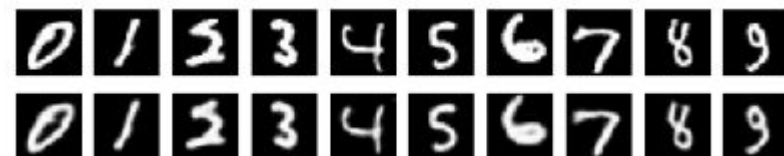
AUTHORS

REIMPLEMENTATION

Accuracy scores

	interpret able	without prototype l.	standard CNN
test set accuracy	99.22%	99.23%	99.24%

Autoencoder samples



	interpret able	without prototype l.	standard CNN
test set accuracy	99.12%	99.37%	99.09%



Learned prototypes



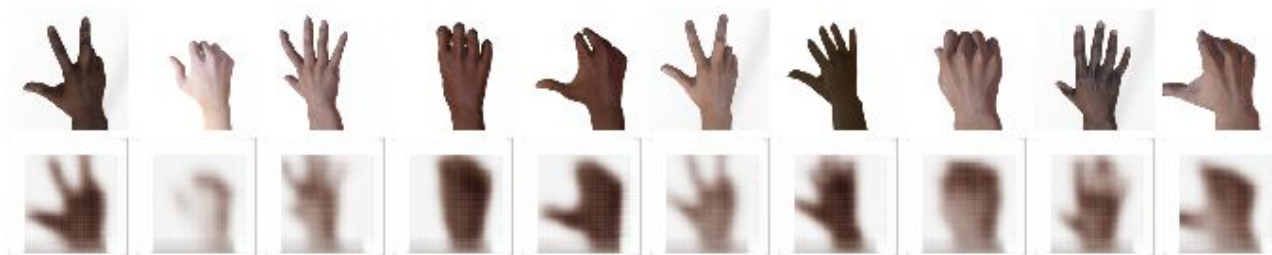
S S P P R

Explainable vs standard classification

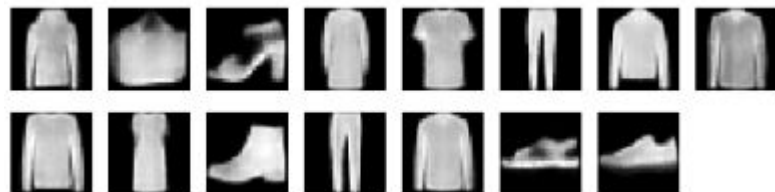
Explainable accuracy: **79.57%**

Standard accuracy: **78.76%**

Autoencoder results

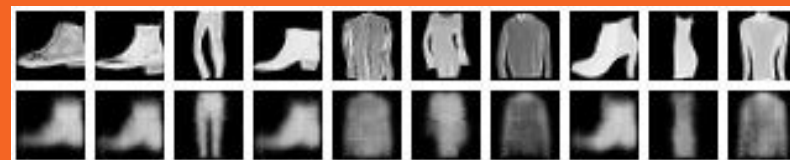
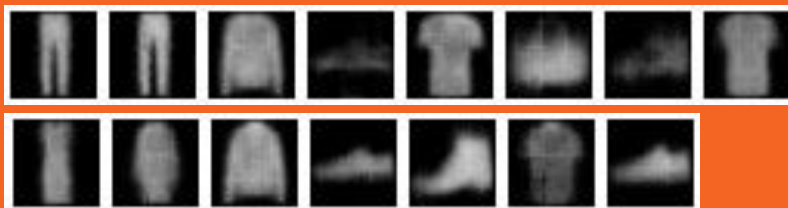


Learned prototypes



Accuracy & Autoencoder samples

	test set accuracy
authors	89.95%
reimplementation	90.23%

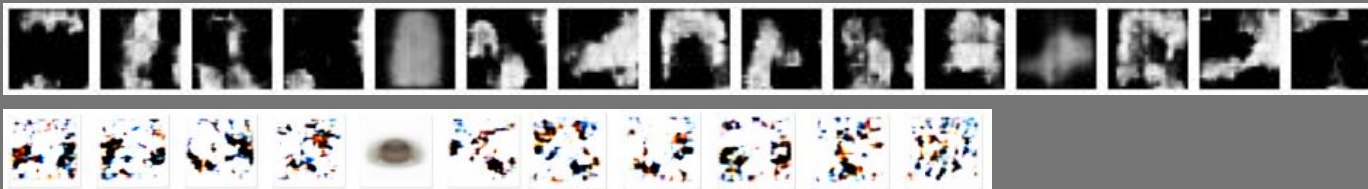


AUTHORS

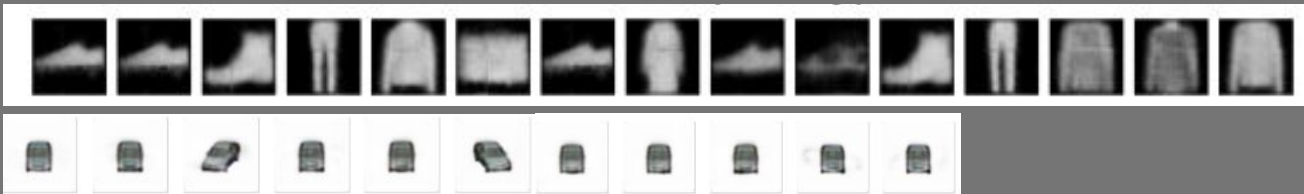
REIMPLEMENTATION

Ablation study

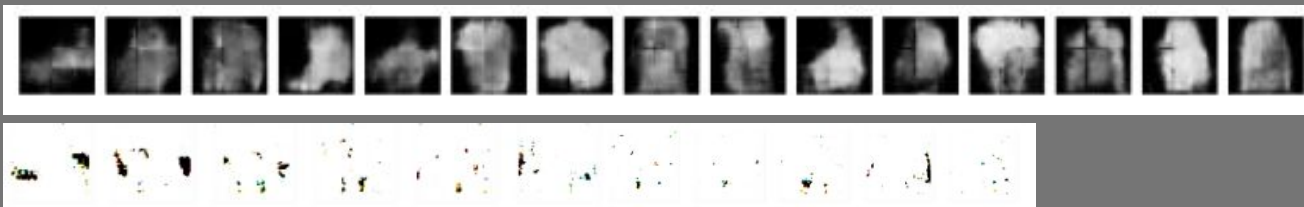
- without R1



- without R2



- without R1 and R2



Any questions?



—

**Thanks
for your
attention!**

