A neuro-inspired computing module for unsupervised continual learning

> Constantine Dovrolis School of Computer Science, Georgia Tech

James Smith



Seth Baer

Zsolt Kira

Supported by the Lifelong Learning Machines (L2M) program of DARPA/MTO

Disclaimers & Apologies

- This work is neuro-inspired but it is not neuromorphic
 - Not at the level of neurons & synapses
- Strong opinion: Neuroscience can provide great insights & ideas to ML/AI beyond the level of neurons and synapses
 - Network neuroscience
 - Cognitive neuroscience
 - Evolutionary neuroscience
 - Neurology & psychiatry

Unsupervised Continual Learning (UCL)

- Learn efficient data representations from stream of unlabeled data
 - No labeled data in stream
 - No specific task or reward signal given
 - Class/data distribution is non-stationary
 - No storage/replay of past inputs

Occasionally, create associations between learned representations and named concepts/classes

Gavenovegenienatiese?abeled data (1 per class?)
Labeled data does not change representations

What UCL is NOT

- * Representation learning
 - Typically does not operate in online context
- * Supervised continual learning
 - Learns continually from stream of labeled data
- Semi-supervised learning
 - Learns from both labeled and unlabeled data
- * Self-supervised learning
 - Relies strictly on prediction does not learn associations

Arguably, UCL is how natural organisms learn most of the time

What we borrow from neuroscience

- Hierarchical organization between cortical areas with feedforward 1. feedback connections
- Mountcastle: cortical columns 2 perform an (unknown) common function across cortex
- The "canonical cortical circuit" of 3. Douglas & Martin
- Intra-column recurrent circuits 4. and role of inhibitory neurons
- Friston: cortical columns perform 5. predictive coding
- Cortical columns appear to learn "prototypes" see Tanaka '96 for monkey IT experiments 6.
- E/I recurrent circuits (similar to 7. L4) can perform k-means online clustering (Pehlevan et al. '18)

Hebbian W

Self-Trained Associative Memory (STAM) Architecture Overview

Illustration of how STAMs work

* Gradual reduction of intrinsic dimensionality

- If there are N centroids at a layer, and that layer consists of M STAMS (or RFs), the output image at that layer can take N^M possible values
- Deeper in hierarchy: N increases, M decreases (so that N^M decreases)

Key equations

* Centroid selection:

$$c(x_{i,m}) = \arg \min_{j=1...|C_i|} ||x_{i,m} - w_{i,j}||$$

* Centroid online learning:

$$w_{i,j} = \alpha x_{i,m} + (1 - \alpha) w_{i,j}$$
, when $c(x_{i,m}) = j$

* Centroid novelty detection:

$$\begin{aligned} \mu_{j} &= \alpha \left| |x_{i,m} - w_{i,j}| \right| + (1 - \alpha) \mu_{j} \\ \hat{\sigma}_{j} &= \alpha \left| \left| |x_{i,m} - w_{i,j}| \right| - \mu_{j} \right| + (1 - \alpha) \hat{\sigma}_{j} \\ &|x_{i,m} - w_{i,j}| \right| > \mu_{j} + 3 \hat{\sigma}_{j} \end{aligned}$$

- * Centroid forgetting:
 - when capacity of N centroids at a layer is reached, forget Least Recently Used centroid

The role of top-down predictions

Phase-1: stream of Os & 1s

Phase 1 Unlabeled Data Stream Samples (5000 Images Seen)

010100100 0011000100

Phase-1 evaluation given a couple of labeled examples

Phase-2: stream of {0,1,2,3}

Phase 2 Unlabeled Data Stream Samples (10000 Images Seen)

Phase-2 evaluation given a couple of labeled examples

Phase-5: stream of {0,1,..9}

Phase 5 Unlabeled Data Stream Samples (45000 Images Seen)

Phase-5 evaluation given a couple of labeled examples per class

Accuracy as the system learns more classes

STAM hierarchy performs best for the entire scenario

Accuracy as function of number of labeled examples

Next Steps

- Work with natural image datasets and video
- * Compare with additional baseline methods compatible with UCL problem
- Generalize classification to use centroids from any layer in hierarchy
- Apply in video and timeseries problems (where prediction and top-down connections can play major role)