

# Task-Driven Dictionary Learning in Distributed Online Settings

Alec Koppel<sup>\*</sup>, Garrett Warnell<sup>†</sup>, Ethan Stump<sup>†</sup>

<sup>\*</sup>University of Pennsylvania, Philadelphia, PA

<sup>†</sup>U.S. Army Research Laboratory, Adelphi, MD

Asilomar Conference on Signals, Systems, and Computers

Pacific Grove, CA, Nov. 10, 2015

- ▶ Goal: visual awareness in mobile robotic teams in unknown domains
  - ⇒ Focus on cases where external state information unavailable
  - ⇒ Little a priori knowledge of environment on platform
- ▶ Online (real-time) training algorithms necessary for this setting
  - ⇒ gain awareness of operating environment
- ▶ Distributed protocol alleviates need for central base station
  - ⇒ better suited to feed into distributed learning-based control

- ▶ Robot  $i$  observes signals  $\theta_{i,t} \in \Theta$ ,  $t = 1, \dots$  based on path it takes
  - ⇒ predict environmental properties  $\mathbf{y}_{i,t} \in \mathcal{Y}$  with this info
  - ⇒ formulate as stoch. opt. problem:  $\min_{\mathbf{x}_i} \mathbb{E}_{\theta_i, \mathbf{y}_i} [f(\mathbf{x}_i; (\theta_i, \mathbf{y}_i))]$
  - ⇒ loss function  $f$ , regressor  $\mathbf{x}_i$  ⇒ discriminative model
- ▶ Wrinkle: individual robots only have info. based on traversed path
  - ⇒ may **omit regions** of feature space crucial for **effective prediction**.
- ▶ Communicate with robotic network
  - ⇒ greater domain "understanding" among individual robots
- ▶ **This work: distributed online predictive algorithms in robotic teams**
- ▶ Develop new capability: individual robots make global inferences
  - ⇒ only observe distinct subsets of feature space

- ▶ If relationship between random pair  $(\theta, \mathbf{y})$  is complicated. . .
  - ⇒ use alternative encoding of  $\theta$  ⇒ reveal latent data structure
- ▶ DSP methods mostly rely on alternative signal representations
  - ⇒ Based on processing task (e.g. Fourier basis, wavelets, PCA)
- ▶ In **dictionary learning**, learn representation directly from data
  - ⇒ **Task-driven**: tailor dictionary to learning task (Mairal '12)
- ▶ We extend task-driven dictionary learning to multi-robot settings
  - ⇒ online perception in mobile robotic teams

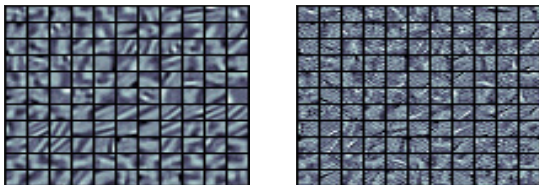


Figure: Initialized (left) and learned (right) dictionary for small image patches.

- ▶ Represent signals  $\theta_t$  as combos. of  $k$  basis elements  $\{\mathbf{d}_l\}_{l=1}^k$ 
  - ⇒ learn dictionary  $\mathbf{D} \in \mathbb{R}^{m \times k}$  from data
  - ⇒ Denote coding (coefficients) of  $\theta_t$  w.r.t. dictionary as  $\alpha_t \in \mathbb{R}^k$
- ▶ **Representation loss**  $g(\alpha_t, \mathbf{D}; \theta_t) \Rightarrow$  small if  $\mathbf{D}\alpha_t$  and  $\theta_t$  close
  - ⇒  $\mathbf{D}\alpha_t$  is representation of  $\theta_t$  w.r.t dictionary  $\mathbf{D}$
- ▶ Formulate the **coding problem** (lasso, elastic net)

$$\alpha^*(\mathbf{D}; \theta_t) := \underset{\alpha_t \in \mathbb{R}^k}{\operatorname{argmin}} g(\alpha_t, \mathbf{D}; \theta_t) .$$

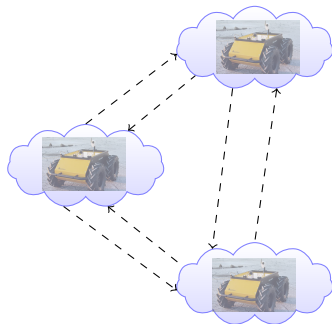
- ▶ Dictionary learning
  - ⇒ seek  $\mathbf{D}$  such that signals  $\theta_t$  well-represented by  $\mathbf{D}\alpha^*(\mathbf{D}; \theta_t)$

- ▶ Tailor dictionary to **discriminative modeling task**
- ▶ Use coding  $\alpha^*(\mathbf{D}; \theta_t)$  as representation of signal  $\theta_t$
- ▶ Decision variable  $\mathbf{x}$   $\Rightarrow$  predict the label/vector  $\mathbf{y}_t$  given  $\alpha^*(\mathbf{D}; \theta_t)$ .
- ▶ **Loss function**  $f(\mathbf{D}, \mathbf{x}; (\theta_t, \mathbf{y}_t)) = f(\alpha^*(\mathbf{D}; \theta_t), \mathbf{D}, \mathbf{x}; (\theta_t, \mathbf{y}_t))$   
 $\Rightarrow$  predictive quality of  $\mathbf{x}$  for output var.  $\mathbf{y}_t$  given coding  $\alpha^*(\mathbf{D}; \theta_t)$
- ▶ Discriminative dictionary learning

$$(\mathbf{D}^*, \mathbf{x}^*) := \operatorname{argmin}_{\mathbf{D} \in \mathcal{D}, \mathbf{x} \in \mathcal{X}} \mathbb{E}_{\theta, \mathbf{y}} \left[ f(\mathbf{D}, \mathbf{x}; (\theta, \mathbf{y})) \right].$$

- $\Rightarrow$  Learn jointly regression weights  $\mathbf{x}$  and dictionary  $\mathbf{D}$
- $\Rightarrow$  Non-convex stochastic program

- ▶ Multi-agent system
  - ⇒ graph  $\mathcal{G} = (V, \mathcal{E})$
  - ⇒  $|V| = N, |\mathcal{E}| = M$
- ▶ Neighborhood of agent  $i$ 
  - ⇒  $n_i = \{j : (j, i) \in \mathcal{E}\}$



- ▶ Each agent  $i$  aims to learn a regressor  $\mathbf{x}$  and dictionary  $\mathbf{D}$ 
  - ⇒ over observations of whole network  $\{(\boldsymbol{\theta}_i, \mathbf{y}_i)\}_{i=1}^N$

$$(\mathbf{x}^*, \mathbf{D}^*) = \operatorname{argmin}_{\mathbf{x}, \mathbf{D}} \sum_{i=1}^N \mathbb{E}_{\boldsymbol{\theta}_i, \mathbf{y}_i} [f(\mathbf{D}, \mathbf{x}; (\boldsymbol{\theta}_i, \mathbf{y}_i))]$$

- ▶ We develop distributed online iterative methods for this problem

- ▶ Incentivize agreement via constraint  $\mathbf{D}_i = \mathbf{D}_j, \mathbf{x}_i = \mathbf{x}_j$  for all  $j \in n_i$
- ▶ Decentralized task-driven dictionary learning problem

$$\{\mathbf{D}_i^*, \mathbf{x}_i^*\}_{i=1}^N := \underset{\mathbf{D}_i \in \mathcal{D}, \mathbf{x}_i \in \mathcal{X}}{\operatorname{argmin}} \sum_{i=1}^N \mathbb{E}_{\boldsymbol{\theta}_i, \mathbf{y}_i} [f(\mathbf{D}_i, \mathbf{x}_i; (\boldsymbol{\theta}_i, \mathbf{y}_i))] .$$

such that  $\mathbf{D}_i = \mathbf{D}_j, \mathbf{x}_i = \mathbf{x}_j$  for all  $j \in n_i$

- ▶ Enforcing agreement constraint would require global coordination  
 $\Rightarrow$  Define stochastic Lagrangian relaxation

$$\hat{\mathcal{L}}_t(\mathbf{D}, \mathbf{x}, \boldsymbol{\Lambda}, \boldsymbol{\nu}) = \sum_{i=1}^N [f(\mathbf{D}_i, \mathbf{x}_i; (\boldsymbol{\theta}_{i,t}, \mathbf{y}_{i,t}))] \\ + \operatorname{tr}(\boldsymbol{\Lambda}^T \mathbf{C}_D \mathbf{D}) + \boldsymbol{\nu}^T \mathbf{C}_x \mathbf{x}$$

- $\Rightarrow$  Apply saddle point to stochastic Lagrangian  $\Rightarrow$  distributed alg.



- ▶ At robot  $i$ , time  $t$ , observe  $(\boldsymbol{\theta}_{i,t}, \mathbf{y}_{i,t})$ ,
- ▶ Compute coding  $\boldsymbol{\alpha}_{i,t+1}^* = \operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^k} g(\boldsymbol{\alpha}, \mathbf{D}_{i,t}; \boldsymbol{\theta}_{i,t})$   
⇒ In practice chosen as *sparse coding* via lasso or elastic-net
- ▶ Update primal variables at robot  $i$

$$\mathbf{D}_{i,t+1} = \mathbf{D}_{i,t} - \epsilon_t \left( \nabla_{\mathbf{D}_i} f_i(\mathbf{D}_{i,t}, \mathbf{x}_{i,t}; (\boldsymbol{\theta}_{i,t}, \mathbf{y}_{i,t})) + \sum_{j \in n_i} (\boldsymbol{\Lambda}_{ij,t} - \boldsymbol{\Lambda}_{ji,t}) \right),$$

$$\mathbf{x}_{i,t+1} = \mathbf{x}_{i,t} - \epsilon_t \left( \nabla_{\mathbf{x}_i} f_i(\mathbf{D}_{i,t}, \mathbf{x}_{i,t}; (\boldsymbol{\theta}_{i,t}, \mathbf{y}_{i,t})) + \sum_{j \in n_i} (\boldsymbol{\nu}_{ij,t} - \boldsymbol{\nu}_{ji,t}) \right),$$

- ▶ Update dual variables at network communication link  $(i, j)$

$$\boldsymbol{\Lambda}_{ij,t+1} = \boldsymbol{\Lambda}_{ij,t} + \epsilon_t (\mathbf{D}_{i,t} - \mathbf{D}_{j,t})$$

$$\boldsymbol{\nu}_{ij,t+1} = \boldsymbol{\nu}_{ij,t} + \epsilon_t (\mathbf{x}_{i,t} - \mathbf{x}_{j,t})$$

- ▶ Decentralized dynamic dictionary learning  $\Rightarrow$  Block saddle point alg.
- ▶ Stochastic approximation:  $\mathcal{L}(\mathbf{D}, \mathbf{w}, \mathbf{\Lambda}, \nu) = \mathbb{E}_{\theta, \mathbf{y}}[\hat{\mathcal{L}}_t(\mathbf{D}, \mathbf{x}, \mathbf{\Lambda}, \nu)]$   
 $\Rightarrow$  primal stochastic gradient descent

$$\mathbf{D}_{t+1} = \mathbf{D}_t - \epsilon_t \nabla_{\mathbf{D}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t),$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \epsilon_t \nabla_{\mathbf{x}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t).$$

$\Rightarrow$  dual stochastic gradient ascent

$$\mathbf{\Lambda}_{t+1} = \mathbf{\Lambda}_t + \epsilon_t \nabla_{\mathbf{\Lambda}} \hat{\mathcal{L}}_t(\mathbf{D}_{t+1}, \mathbf{x}_{t+1}, \mathbf{\Lambda}_t, \nu_t),$$

$$\nu_{t+1} = \nu_t + \epsilon_t \nabla_{\nu} \hat{\mathcal{L}}_t(\mathbf{D}_{t+1}, \mathbf{x}_{t+1}, \mathbf{\Lambda}_t, \nu_t).$$

- ▶  $\nabla_{\mathbf{D}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t) \Rightarrow$  Projected stoch. Lagrangian grad. w.r.t.  $\mathbf{D}$   
 $\Rightarrow$  gradient approximated with current signals  $\{\theta_{i,t}, \mathbf{y}_{i,t}\}_{i=1}^N$

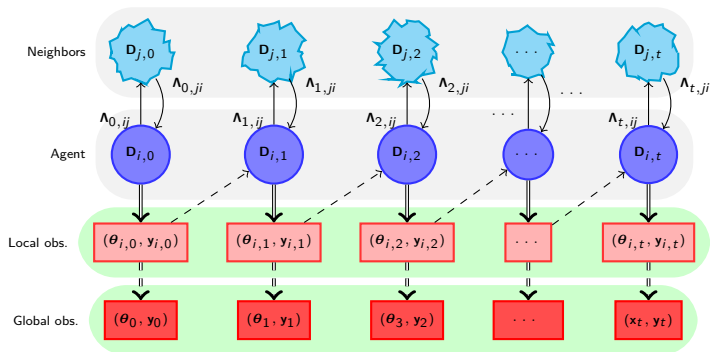
- ▶ Network  $\mathcal{G} \Rightarrow$  symmetric and connected with diameter  $D$ .
- ▶ Diminishing step-size rules:  $\sum_{t=0}^{\infty} \epsilon_t = \infty$  and  $\sum_{t=0}^{\infty} \epsilon_t^2 < \infty$
- ▶ Mean and variance conditions of Lagrangian stochastic gradients

$$\mathbb{E}[\|\delta_{\mathbf{D},t}\| \mid \mathcal{F}_t] \leq A\epsilon_t,$$

$$\mathbb{E}[\|\nabla_{\mathbf{D}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{x}_t, \boldsymbol{\Lambda}_t, \nu_t)\|^2 \mid \mathcal{F}_t] \leq \sigma^2.$$

- ▶ Feasible dictionary set is restricted to those with unit column-norms

$$\mathcal{D} = \{\mathbf{D} \in \mathbb{R}^{m \times k} : \|\mathbf{d}_j\| \leq 1, j = 1 \dots k\}.$$



- ▶ Dictionary learning scheme depicted above  
 ⇒ model parameters work in same manner
- ▶ only exchange local decision variables and Lagrange multipliers

## Theorem

Saddle pt. seq.  $(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t)$  *converges to stationarity* in expectation:

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\mathbf{D}} \mathcal{L}(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t)\|] = 0,$$
$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t)\|] = 0$$

*Asymptotic feasibility condition* achieved in expectation:

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\mathbf{\Lambda}} \mathcal{L}(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t)\|] = 0$$
$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\nu} \mathcal{L}(\mathbf{D}_t, \mathbf{x}_t, \mathbf{\Lambda}_t, \nu_t)\|] = 0$$

- ▶ Performance guarantee for D4L
  - ⇒ convergence in non-convex stochastic opt.
  - ⇒ sensitive to data distribution, step-size, network structure

- ▶ Texture database classification problem  $\Rightarrow$  Brodatz textures
  - $\Rightarrow$  Insight into **dynamic image processing problems**
  - $\Rightarrow$  **Toy model of real-time navigability analysis** in robotic teams
- ▶ Real-time image data  $\Rightarrow$  train multi-class logistic regression weights
- ▶ Decentralized dynamic texture classification
  - $\Rightarrow$  Subset of textures: {grass, bark, straw, herringbone\_weave}

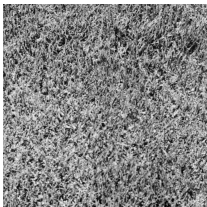


Figure: Sample images from Brodatz textures.

- ▶ Multi-class logistic regression prob.  $\Rightarrow$  Robot  $i$  receives signals  $\theta_{i,t}$   
 $\Rightarrow$  output a **decision variable**  $\mathbf{y}_{i,t} \in \{0, 1\}^C \Rightarrow C$  no. of classes
- ▶  $[\mathbf{y}_{i,t}]_c \Rightarrow$  binary indicator of whether signal falls in class  $c$ .
- ▶ Local loss  $f_i \Rightarrow$  negative log-likelihood of prob. model

$$f_i(\mathbf{D}_i, \mathbf{X}_i; (\theta_i, \mathbf{y}_i)) = \log \left( \sum_{c=1}^C e^{\mathbf{x}_{i,c}^T \boldsymbol{\alpha}_i^* + x_{i,c}^0} \right) - \sum_{c=1}^C \left( y_{i,c} \mathbf{x}_{i,c}^T \boldsymbol{\alpha}_i^* + w_{i,c}^0 \right) + \xi \|\mathbf{X}_i\|_F^2,$$

- ▶  $\boldsymbol{\alpha}_i^* \Rightarrow$  sparse coding via elastic-net min. prob.
- ▶  $g_c(\boldsymbol{\alpha}_i^*) = e^{\mathbf{x}_{i,c}^T \boldsymbol{\alpha}_i^* + x_{i,c}^0}$  is activation function;  
 $\Rightarrow g_c(\mathbf{z}_i) / \sum_{c'} g_{c'}(\mathbf{z}_i) \Rightarrow$  prob.  $\mathbf{z}_i$  in class  $c$   
 $\Rightarrow \mathbf{z}_i \Rightarrow$  average of image sub-patches
- ▶ Classification decision  $\Rightarrow$  **maximum likelihood class label**  
 $\Rightarrow \tilde{c} = \operatorname{argmax}_c g_c(\mathbf{z}_i) / \sum_{c'} g_{c'}(\mathbf{z}_i); \quad [\mathbf{y}_{i,t}]_c = 0$  for  $c \neq \tilde{c}$

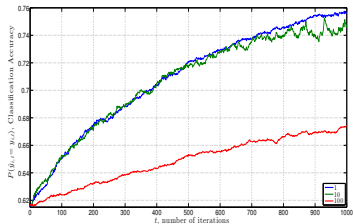
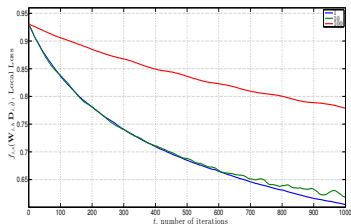


Figure: Local loss (left) and classification accuracy (right) versus iteration  $t$ .

- ▶ Slower learning in larger networks ( $N \uparrow$ )
- ▶ Non-convexity hurts more in larger-networks
  - ⇒ Smaller step-sizes required for convergence



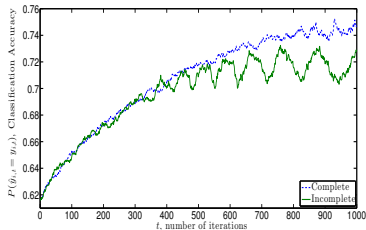
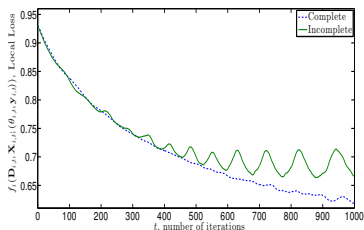


Figure: Log-likelihood (left) and classification accuracy (right) vs. time  $t$ .

- ▶  $N = 10$  node random network, results shown for random  $j \in V$
- ▶ Agents observe **random incomplete subsets** of feature space
- ▶ Still learn **global information** and reach consensus
- ▶ Moderate classifier performance
  - ⇒ due to small step-size required for convergence
  - ⇒ Small step-sizes required for convergence

- ▶  $N = 3$  network of Clearpath Huskies equipped with Zigby radios
  - ⇒ observes images and pose information online via on-board sensors
  - ⇒ each robot partition them into small patches ⇒ classify patches.
- ▶ True labels generated by  $\mathbf{1}\{|\omega - \hat{\omega}| > \gamma\}$ 
  - ⇒ Threshold difference of commanded and measured state info.
  - ⇒ Decentralized maximum likelihood prediction of robot slipping
- ▶ Experiments at Lejeune Robotics Test Facility ⇒ Thanks to ARL!

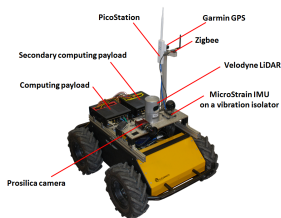


Figure: Sample image (left) from a  $N = 3$  robot network of Huskies (right).

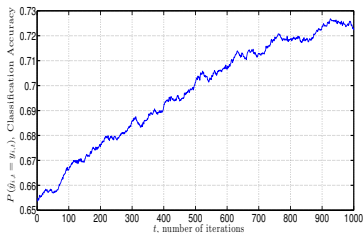
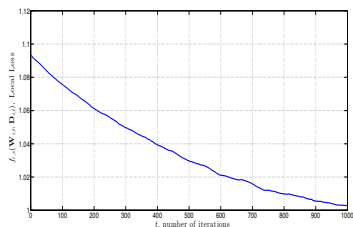


Figure: Log-likelihood (left) and classification accuracy (right) versus time  $t$ .

- ▶ Experimental setting:  $N = 3$  complete graph
- ▶ Label classes: grass and asphalt
- ▶ Robotic implementation  $\Rightarrow$  promising initial results

- ▶ Pattern recognition  $\Rightarrow$  finding good signal representation
- ▶ Online task-driven dictionaries  $\Rightarrow$  visual awareness in robotic teams
- ▶ Decentralized non-convex stochastic opt. problem
- ▶ **Block-variant of saddle pt. method  $\Rightarrow$  convergence in expectation**
- ▶ Implementation on robotic network of Huskies
  - $\Rightarrow$  distributed online protocol for gaining environmental awareness

<http://seas.upenn.edu/~akoppel/>

- ▶ A. Koppel, G. Warnell, E. Stump, and A. Ribeiro, “D4L: Decentralized Dynamic Discriminative Dictionary Learning,” in Proc. Int. Conf. Intelligent Robotics and Systems, Hamburg, Germany, Sep 28-Oct2 2015
- ▶ A. Koppel, G. Warnell, and E. Stump. “A Stochastic Primal-Dual Algorithm for Task-Driven Dictionary Learning in Networks.” in Proc. Asilomar Conf. on Signals Systems Computers, Pacific Grove, CA, November 8-11 2015.
- ▶ A. Koppel, G. Warnell, E. Stump, and A. Ribeiro, “D4L: Decentralized Dynamic Discriminative Dictionary Learning,” IEEE Trans. Signal Process., July. 2015. (submitted).