Slow Feature Analysis For Curiosity-Driven Agents

Varun Raj Kompella (www.idsia.ch/~kompella)

Co-organizers: Matthew Luciw and Juergen Schmidhuber

The Swiss AI Lab IDSIA, Switzerland



Lugano, Switzerland





http://www.atlanticolugano.ch/images/lugano05.jpg

Co-Organizers Dr. Matthew D. Luciw

□ Webpage: <u>www.idsia.ch/~luciw</u>

- Currently working at Neurala, Inc, Boston US.
- Offshoot of Boston University neuromorphic research
- Working with NASA, USAF and private customers to design deep learning autonomy algorithms for ground/air robots.







www.neurala.com

Co-Organizers Prof. Juergen Schmidhuber

- Director of the Swiss AI Lab IDSIA, Lugano, Switzerland.
- Professor of AI at the Faculty of Computer Science at the University of Lugano (USI) and SUPSI.
- He is known for his work on machine learning, artificial intelligence, artificial neural networks, digital physics and low-complexity art.
 - Webpage: <u>www.idsia.ch/~juergen</u>





Tutorial Webpage

- <u>http://www.idsia.ch/~kompella/mywork/WCCI_Tutorial.html</u>
- Related material: Slides, Papers, Videos, Codes, Online-browser implementations.
- □ Will be constantly updated for new material.

Motivation & Background









Sneak Preview (What to Expect)



Sneak Preview (What to Expect)



Anthropomorphic Robots (ASIMO, PETMAN, iCUB)







 Reinforcement Learning -Task Specific Behaviors via trial-error interactions



- Reinforcement Learning -Task Specific Behaviors via trial-error interactions
- Difficulty: Pure randomexploration in large sensory and joint state spaces.
 - □ eg: (320 x 240) x (~50) continuous state variables!



- Reinforcement Learning -Task Specific Behaviors via trial-error interactions
 - Difficulty: Pure randomexploration in large sensory and joint state spaces.
 - eg: (320 x 240) x (~50)
 continuous state variables!
 - Self supervised learning







Wait.... Is this Interesting ??





= What is this ??

How can I push the cup? grasp the cup?

Supervised Learning (SL)

Training examples: $(x_i, y_i), x_i \sim D(x), y_i = f(x_i)$. x-samples, y-labels.

Learning Objective: Improve a hypothesis h(x) by minimising a loss function L(y, h(x)).

Examples: Artificial Neural Networks, Support Vector Machine



Unsupervised Learning (UL)

Training examples: (x_i) , $x_i \sim D(x)$. x-samples.

Learning Objective: Minimise a loss function that reflects some statistical structure of the input.

Examples: Principal Component Analysis, Hidden Markov Models



Reinforcement Learning (RL)



- Learning Objective: Take actions to maximise expected cumulative rewards.
- Value Function: A number that denotes how valuable a state-action is in solving the objective.
- □ **Policy**: Mapping between states and actions.
- **Types of RL**: Batch, Open-ended.

Markov Decision Process (MDP)

Defined as a 4-tuple (S, A, P, R)
S - Fully observable state space; A - Space of actions
P - Stationary state transition fn. ; R - Stationary reward fn.
Possesses the Markov Property



Intrinsic Motivation





Provides a mathematical formalism for describing curiosity and creativity.

- Provides a mathematical formalism for describing curiosity and creativity.
- □ Interesting: Learnable but as-yet-unknown aspects of the environment.

- Provides a mathematical formalism for describing curiosity and creativity.
- □ Interesting: Learnable but as-yet-unknown aspects of the environment.
- □ **Un-interesting**: Predictable or inherently unlearnable (noise)

- Provides a mathematical formalism for describing curiosity and creativity.
- □ Interesting: Learnable but as-yet-unknown aspects of the environment.
- Un-interesting: Predictable or inherently unlearnable (noise)
- Curiosity Rewards: Proportional to the improvement of an internal model/predictor of the environment.
















Artificial Curiosity

Now this slide is boring! Move on!



CDRL objective : Maximise accumulation of expected cumulative curiosity rewards.



- CDRL objective : Maximise accumulation of expected cumulative curiosity rewards.
 - □ Agent is motivated to explore wherever it is makes maximum learning progress.

- CDRL objective : Maximise accumulation of expected cumulative curiosity rewards.
 - ☐ Agent is motivated to explore wherever it is makes maximum learning progress.
- **Task-independent** intrinsic reward function

- CDRL objective : Maximise accumulation of expected cumulative curiosity rewards.
 - ☐ Agent is motivated to explore wherever it is makes maximum learning progress.
- **Task-independent** intrinsic reward function
- Challenges in CDRL: High dimensionality of the state spaces, nonstationary rewards

High-Dimensionality of State-Space



640 x 480 Dim

High-Dimensionality of State-Space



640 x 480 Dim

Principal Component Analysis (PCA)

Unsupervised Learning method

- Extracts linear uncorrelated feature vectors, which represent directions of maximal variance in the input data.
- Eigen-Decomposition
- Dimensionality reduction:
 Project input data on fewer
 PCA vectors (m < n):

 $y = V_{pca}^{T} \cdot x$ (m x 1) (m x n) (n x 1)





Slow Feature Analysis (Wiskott and Sejnowski, 2002)



- Unsupervised learning method
- Uses temporal coherence in the input data
- Extracts slowly varying components from rapidly changing raw sensory inputs (Slowness Learning)



SFA Formalised

Given an I-dimensional input signal $\mathbf{x}(t)$, find a set of J instantaneous real-valued functions $\mathbf{g}(x)$, which generate a J-dimensional output signal $\mathbf{y}(t)$, such that for each $j \in \{1, ..., J\}$

Definition

$$\Delta_j := \Delta(y_j) := \langle \dot{y}_j^2 \rangle$$
 is minimal, with the constraints

$$\langle y_j \rangle = 0$$
 (zero mean), (1)
 $\langle y_j^2 \rangle = 1$ (unit variance), (2)
 $\langle i < j : \langle y_i y_j \rangle = 0$ (decorrelation and order), (3)

❑ Whitening: Input is normalised to have zero mean and unit variance via PCA-1 (V₁)

- Whitening: Input is normalised to have zero mean and unit variance via PCA-1 (V1)
- Derivative: A fast approximation of the whitened input is computed using backward-difference.

- Whitening: Input is normalised to have zero mean and unit variance via PCA-1 (V1)
- Derivative: A fast approximation of the whitened input is computed using backward-difference.
- Minor Components: Extract components with least Eigen-values of the derivative signal using PCA-2 (V₂)

- Whitening: Input is normalised to have zero mean and unit variance via PCA-1 (V1)
- Derivative: A fast approximation of the whitened input is computed using backward-difference.
- Minor Components: Extract components with least Eigen-values of the derivative signal using PCA-2 (V₂)
- **Slow Features:** $V_{sf} = V_1 \bullet V_2$

- Whitening: Input is normalised to have zero mean and unit variance via PCA-1 (V1)
- Derivative: A fast approximation of the whitened input is computed using backward-difference.
- Minor Components: Extract components with least Eigen-values of the derivative signal using PCA-2 (V₂)
- **Slow Features:** $V_{sf} = V_1 \bullet V_2$
- **SFA Outputs**: $y(t) = V_{sf}^{T} \bullet x(t)$

SFA Examples



$$\breve{x}_1(t) = \sin(t) + \cos(11 t)^2,$$

 $\breve{x}_2(t) = \cos(11 t), t \in [0, 2\pi],$



SFA Examples

Environment

Input Image Stream



First Slow Feature Output

SFA Examples





□ Laplacian Eigen Maps (LEMs):

□ Laplacian Eigen Maps (LEMs):

 (a) A low-dimensional embedding of the state-space, which captures the local-connectedness between the states (Adjacency).



- □ Laplacian Eigen Maps (LEMs):
 - (a) A low-dimensional embedding of the state-space, which captures the local-connectedness between the states (Adjacency).
 - (b) Useful for representing valuefunctions for RL (Proto-value functions).



- Laplacian Eigen Maps (LEMs):
 - (a) A low-dimensional embedding of the state-space, which captures the local-connectedness between the states (Adjacency).
 - (b) Useful for representing valuefunctions for RL (Proto-value functions).
 - (c) LEMs are non-parametric dimensionality scales with the number of data points.



 SFA is a low-complex linear function approximation of LEMs





- SFA is a low-complex linear function approximation of LEMs
- SFA is a parametric method that generalises to unseen data.





- SFA is a low-complex linear function approximation of LEMs
- SFA is a parametric method that generalises to unseen data.
- Can therefore be a good preprocessor for RL (Legenstein et al. 2010, Kompella et al. 2011, Luciw and Schmidhuber, 2012).





SFA for Open-Ended Online RL?



SFA for Open-Ended Online RL?

- □ Limitations of SFA for open-ended online learning RL agents:
 - □ Estimates covariance matrices via batch processing.
 - □ Cubic computational complexity.

Incremental Slow Feature Analysis (IncSFA)



□ Incremental SFA (IncSFA) is an online implementation of SFA

□ Incremental SFA (IncSFA) is an online implementation of SFA

□ Adaptive to changing input statistics

□ Incremental SFA (IncSFA) is an online implementation of SFA

- Adaptive to changing input statistics
- □ Linear computational complexity

- □ Incremental SFA (IncSFA) is an online implementation of SFA
 - Adaptive to changing input statistics
 - □ Linear computational complexity
 - Reduced sensitivity to outliers
Incremental Slow Feature Analysis (Kompella et al., 2011)

- □ Incremental SFA (IncSFA) is an online implementation of SFA
 - Adaptive to changing input statistics
 - □ Linear computational complexity
 - □ Reduced sensitivity to outliers
 - □ Adds to the biological plausibility of batch SFA

Batch SFA Re-Visited

- □ Whitening: Input is normalised to have zero mean and unit variance via PCA-1 (V₁)
- □ **Minor Components**: Extract components with least Eigen-values of the derivative signal using PCA-2 (V₂)
- $\Box \text{ Slow Features: } V_{sf} = V_1 \bullet V_2$

IncSFA Simplified



IncSFA Simplified

Incremental whitening using Candid Covariance-free Incremental PCA (CCIPCA).



IncSFA Simplified

- Incremental whitening using Candid Covariance-free Incremental PCA (CCIPCA).
- Incremental extraction of minor components using Covariance-free Incremental MCA (CIMCA).

Candid Covariance-free Incremental PCA (CCIPCA; Weng et al.)

 Updates eigenvalues and eigenvectors from each centred observation.

 Algorithm 2: CCIPCA-Update(V, K, u, η)

 //Candid Covariance-Free Incremental PCA

 1 $u_1 \leftarrow u$

 2 for $i \leftarrow l$ to K do

 3
 //Principal component update

 3
 //Principal component update

 3
 //Residual

 4
 $u_i \leftarrow (1 - \eta) v_i + \eta \left[\frac{u_i \cdot v_i}{\|v_i\|} u_i \right]$

 //Residual

 4
 $u_{i+1} = u_i - \left(u_i^T \frac{v_i}{\|v_i\|} \right) \frac{v_i}{\|v_i\|}$

 5 end

 6 return V

Candid Covariance-free Incremental PCA (CCIPCA; Weng et al.)

Updates eigenvalues and eigenvectors from each centred observation.

Hebbian Update ($\Delta w \propto \eta x_i y_i$) + Residual Algorithm 2: CCIPCA-Update(V, K, u, η)//Candid Covariance-Free Incremental PCA1 $u_1 \leftarrow u$ 2 for $i \leftarrow l$ to K do3//Principal component update3 $v_i \leftarrow (1 - \eta) v_i + \eta \left[\frac{u_i \cdot v_i}{||v_i||} u_i \right]$ 4 $u_{i+1} = u_i - \left(u_i^T \frac{v_i}{||v_i||} \right) \frac{v_i}{||v_i||}$ 5 end6 return V

Candid Covariance-free Incremental PCA (CCIPCA; Weng et al.)

- Updates eigenvalues and eigenvectors from each centred observation.
- □ Hebbian Update ($\Delta w \propto \eta x_i y_i$) + Residual
- □ Intuitive to tune learning rate

Algorithm 2: CCIPCA-Update(V, K, u, η)//Candid Covariance-Free Incremental PCA1 $u_1 \leftarrow u$ 2 for $i \leftarrow I$ to K do3//Principal component update3 $v_i \leftarrow (1 - \eta) v_i + \eta \left[\frac{u_i \cdot v_i}{\|v_i\|} u_i \right]$ 4 $u_{i+1} = u_i - \left(u_i^T \frac{v_i}{\|v_i\|} \right) \frac{v_i}{\|v_i\|}$ 5 end6 return V

Covariance-free Minor Component Analysis (CIMCA; Peng et al.)

Based on Peng's MCA modifications:
 (a) Covariance free
 (b) Normalisation step

Algorithm 3: CIMCA-Update($\hat{\phi}, J, \dot{z}, \gamma, \eta$)//Covariance-Free Incremental MCA// $\hat{\phi}_i$: i^{th} column of $\hat{\phi}$ 1 $\mathbf{l}_1 \leftarrow 0$ 2 for $i \leftarrow 1$ to J do//Minor component update3 $\hat{\phi}_i \leftarrow (1 - \eta)\hat{\phi}_i - \eta \left[(\dot{z} \cdot \hat{\phi}_i) \dot{z} + \mathbf{l}_i \right].$ //Normalize4 $\hat{\phi}_i \leftarrow \hat{\phi}_i / || \hat{\phi}_i ||.$ //Lateral competition from "lower" components5 $\mathbf{l}_{i+1} \leftarrow \gamma \sum_{j}^{i} (\hat{\phi}_j \cdot \hat{\phi}_i) \hat{\phi}_j$ 6 end7 return W

Covariance-free Minor Component Analysis (CIMCA; Peng et al.)

- Based on Peng's MCA modifications:

 (a) Covariance free
 (b) Normalisation step
- Anti-Hebbian + Sequential addition.

Algorithm 3: CIMCA-Update($\hat{\phi}$, J, \dot{z} , γ , η) //Covariance-Free Incremental MCA $//\hat{\phi_i}$: i^{th} column of $\hat{\phi}$ $1 \mathbf{l}_1 \leftarrow 0$ 2 for $i \leftarrow l$ to J do //Minor component update $\hat{\phi}_i \leftarrow (1-\eta)\hat{\phi}_i - \eta \left[(\dot{\mathbf{z}} \cdot \hat{\phi}_i) \, \dot{\mathbf{z}} + \mathbf{l}_i \right].$ 3 //Normalize $\hat{\phi}_i \leftarrow \hat{\phi}_i / \|\hat{\phi}_i\|.$ 4 //Lateral competition from "lower" components $\mathbf{l}_{i+1} \leftarrow \gamma \, \sum_{j=1}^{i} (\hat{\phi}_j \cdot \hat{\phi}_i) \hat{\phi}_j$ 5 6 end 7 return W

Covariance-free Minor Component Analysis (CIMCA; Peng et al.)

- Based on Peng's MCA modifications:

 (a) Covariance free
 (b) Normalisation step
- Anti-Hebbian + Sequential addition.
- Constant learning rate no need for dynamic rate scheduling.

Algorithm 3: CIMCA-Update($\hat{\phi}$, J, \dot{z} , γ , η) //Covariance-Free Incremental MCA $//\hat{\phi_i}$: i^{th} column of $\hat{\phi}$ $\mathbf{1} \mathbf{l}_1 \leftarrow \mathbf{0}$ 2 for $i \leftarrow l$ to J do //Minor component update $\hat{\phi}_i \leftarrow (1-\eta)\hat{\phi}_i - \eta \left[(\dot{\mathbf{z}} \cdot \hat{\phi}_i) \, \dot{\mathbf{z}} + \mathbf{l}_i \right].$ 3 //Normalize $\hat{\phi}_i \leftarrow \hat{\phi}_i / \|\hat{\phi}_i\|.$ 4 //Lateral competition from "lower" components $\mathbf{l}_{i+1} \leftarrow \gamma \, \sum_{j=1}^{i} (\hat{\phi}_j \cdot \hat{\phi}_i) \hat{\phi}_j$ 5 6 end 7 return W

//Incremental update of J slow features from samples $\mathbf{x} \in \mathscr{R}^I$

- //V : K columns: PCs of x
- $//\hat{\phi}$: J columns: SFs
- $//v^{\gamma}$: First PC in \dot{z} -space
- $//\bar{x}$: Mean of x
- 1 {V, $\hat{\phi}$, \mathbf{v}^{γ} , $\mathbf{\bar{x}}$ } \leftarrow Initialize ()

//Incremental update of J slow features from samples $\mathbf{x} \in \mathscr{R}^I$

- //V : K columns: PCs of x
- $//\hat{\phi}$: J columns: SFs
- $//\boldsymbol{v}^{\boldsymbol{\gamma}}$: First PC in $\dot{z}\text{-space}$
- $//\bar{x}$: Mean of x
- 1 {V, $\hat{\phi}$, \mathbf{v}^{γ} , $\mathbf{\bar{x}}$ } \leftarrow Initialize ()
- 2 for $t \leftarrow 1$ to ∞ do

- 3 $\mathbf{x} \leftarrow \text{Sense}(worldstate)$
- 4 { $\eta_t^{PCA}, \eta_t^{MCA}$ } \leftarrow LrnRateSchedule (θ, t)

Sense

 $//Incremental update of J slow features from samples <math>\mathbf{x} \in \mathscr{R}^{I}$ $//V : K columns: PCs of \mathbf{x}$ $//\hat{\phi} : J columns: SFs$ $//\mathbf{v}^{\gamma} : First PC in \dot{\mathbf{z}} - space$ $//\bar{\mathbf{x}} : Mean of \mathbf{x}$ $1 \{\mathbf{V}, \hat{\phi}, \mathbf{v}^{\gamma}, \bar{\mathbf{x}}\} \leftarrow Initialize ()$ $2 \text{ for } t \leftarrow I \text{ to } \infty \text{ do}$ $3 \mid \mathbf{x} \leftarrow Sense(worldstate)$ $4 \mid \{\eta_{t}^{PCA}, \eta_{t}^{MCA}\} \leftarrow LrnRateSchedule (\theta, t)$

- 5 $\bar{\mathbf{x}} \leftarrow (1 \eta_t^{PCA}) \, \bar{\mathbf{x}} + \eta_t^{PCA} \, \mathbf{x} / / \text{Update mean}$
- 6 $\mathbf{u} \leftarrow (\mathbf{x} \bar{\mathbf{x}}) / / Centering$

Sense

Update mean & de-mean the input

//Incremental update of J slow features from samples $\mathbf{x} \in \mathscr{R}^I$ //V : K columns: PCs of x $//\hat{\phi}$: J columns: SFs $//v^{\gamma}$: First PC in \dot{z} -space $//\bar{x}$: Mean of x 1 {V, $\hat{\phi}$, \mathbf{v}^{γ} , $\mathbf{\bar{x}}$ } \leftarrow Initialize () 2 for $t \leftarrow 1$ to ∞ do $\mathbf{x} \leftarrow \text{Sense}(worldstate)$ 3 $\{\eta_t^{PCA}, \eta_t^{MCA}\} \leftarrow \text{LrnRateSchedule}(\theta, t)$ 4 $\mathbf{\bar{x}} \leftarrow (1 - \eta_t^{PCA}) \, \mathbf{\bar{x}} + \eta_t^{PCA} \, \mathbf{x} \, / \, / \, \text{Update mean}$ 5 $\mathbf{u} \leftarrow (\mathbf{x} - \bar{\mathbf{x}}) / / Centering$ 6 //Candid Covariance-Free Incremental PCA $\mathbf{V} \leftarrow \text{CCIPCA-Update} (\mathbf{V}, K, \mathbf{u}, \eta_t^{PCA})$ 7 S ← ConstructWhiteningMatrix (V) 8 If t > 1 then $(\mathbf{z}_{prev} \leftarrow \mathbf{z}_{curr})$ //Store prev. 9 //Whitening and dim. reduction $\mathbf{z}_{curr} \leftarrow \mathbf{S}^T \mathbf{u}$ 10

Sense

Update mean & de-mean the input

Update whitening matrix & normalize the input

//Incremental update of J slow features from samples $\mathbf{x} \in \mathscr{R}^{I}$ //V : K columns: PCs of x $//\hat{\phi}$: J columns: SFs $//v^{\gamma}$: First PC in \dot{z} -space $//\bar{x}$: Mean of x 1 { $\mathbf{V}, \hat{\phi}, \mathbf{v}^{\gamma}, \bar{\mathbf{x}}$ } \leftarrow Initialize () 2 for $t \leftarrow 1$ to ∞ do $\mathbf{x} \leftarrow \text{Sense}(worldstate)$ 3 $\{\eta_t^{PCA}, \eta_t^{MCA}\} \leftarrow \text{LrnRateSchedule}(\theta, t)$ 4 $\bar{\mathbf{x}} \leftarrow (1 - \eta_t^{PCA}) \, \bar{\mathbf{x}} + \eta_t^{PCA} \, \mathbf{x} \, / \, / \, \text{Update mean}$ 5 $\mathbf{u} \leftarrow (\mathbf{x} - \bar{\mathbf{x}}) / / Centering$ 6 //Candid Covariance-Free Incremental PCA $\mathbf{V} \leftarrow \text{CCIPCA-Update}(\mathbf{V}, K, \mathbf{u}, \eta_t^{PCA})$ 7 $S \leftarrow ConstructWhiteningMatrix (V)$ 8 If t > 1 then $(\mathbf{z}_{prev} \leftarrow \mathbf{z}_{curr})$ //Store prev. 9 //Whitening and dim. reduction $\mathbf{z}_{curr} \leftarrow \mathbf{S}^T \mathbf{u}$ 10 if t > 1 then 11 $\mathbf{z} \leftarrow \left(\mathbf{z}_{curr} - \mathbf{z}_{prev}\right)$ //Approx. derivative 12

Sense

Update mean & de-mean the input

Update whitening matrix & normalize the input

Compute derivative of whitened input

//Incremental update of J slow features from samples $\mathbf{x} \in \mathscr{R}^I$ //V : K columns: PCs of x $//\hat{\phi}$: J columns: SFs $//v^{\gamma}$: First PC in \dot{z} -space $//\bar{x}$: Mean of x 1 { $\mathbf{V}, \hat{\phi}, \mathbf{v}^{\gamma}, \bar{\mathbf{x}}$ } \leftarrow Initialize () 2 for $t \leftarrow 1$ to ∞ do $\mathbf{x} \leftarrow \text{Sense}(worldstate)$ 3 $\{\eta_t^{PCA}, \eta_t^{MCA}\} \leftarrow \text{LrnRateSchedule}(\theta, t)$ 4 $\bar{\mathbf{x}} \leftarrow (1 - \eta_t^{PCA}) \bar{\mathbf{x}} + \eta_t^{PCA} \mathbf{x} / / \text{Update mean}$ 5 $\mathbf{u} \leftarrow (\mathbf{x} - \bar{\mathbf{x}}) / / Centering$ 6 //Candid Covariance-Free Incremental PCA $\mathbf{V} \leftarrow \text{CCIPCA-Update}(\mathbf{V}, K, \mathbf{u}, \eta_t^{PCA})$ 7 S ← ConstructWhiteningMatrix (V) 8 If t > 1 then $(\mathbf{z}_{prev} \leftarrow \mathbf{z}_{curr})$ //Store prev. 9 //Whitening and dim. reduction $\mathbf{z}_{curr} \leftarrow \mathbf{S}^T \mathbf{u}$ 10 if t > 1 then 11 $\mathbf{z} \leftarrow (\mathbf{z}_{curr} - \mathbf{z}_{prev}) / / Approx.$ derivative 12 //For seq. addition (γ) $\mathbf{v}^{\gamma} \leftarrow \text{CCIPCA-Update} (\mathbf{v}^{\gamma}, 1, \mathbf{z}, \eta_t^{PCA})$ 13 $\gamma \leftarrow \mathbf{v}^{\gamma} / \| \mathbf{v}^{\gamma} \|$ 14 //Covariance-free Incremental MCA $\hat{\phi} \leftarrow \text{CIMCA-Update}(\hat{\phi}, J, \dot{z}, \gamma, \eta_t^{MCA})$ 15 end 16 $\mathbf{y} \leftarrow \mathbf{z}_{curr}^T \hat{\boldsymbol{\phi}}$ //Slow feature output 17 18 end

Sense

Update mean & de-mean the input

Update whitening matrix & normalize the input

Compute derivative of whitened input

Update minor component vectors

Experimental Results



http://www.idsia.ch/~kompella/mywork/WCCI_Tutorial.html

Experimental Results



http://www.idsia.ch/~kompella/mywork/WCCI_Tutorial.html

Take-Home Message & Questions?



Take-Home Message & Questions?

□ IncSFA is a low-complex incremental slow feature analysis algorithm.

Take-Home Message & Questions?

- □ IncSFA is a low-complex incremental slow feature analysis algorithm.
- \Box FAQs:
 - □ Is IncSFA truly incremental?
 - □ How easy is it to set the learning rates?
 - Do IncSFA features always converge to their batch counterparts?
 - □ Is IncSFA code available for public?
 - **Python**: <u>http://www.idsia.ch/~kompella/codes/IncSFA.zip</u>
 - □ Matlab: <u>http://www.idsia.ch/~luciw/incsfa.html</u>

5 minutes break!





Recap

Task: An *online* learning framework that is

- curiosity-driven and
- enables skills acquisition from
- high-dimensional video inputs for humanoid robots.

Recap

Task: An *online* learning framework that is

- curiosity-driven and
- enables skills acquisition from
- high-dimensional video inputs for humanoid robots.

How to make the robot curious?

- The theory of Artificial Curiosity

Recap

Task: An *online* learning framework that is

- curiosity-driven and
- enables skills acquisition from
- high-dimensional video inputs for humanoid robots.

How to make the robot curious?

- The theory of Artificial Curiosity

□ How to handle high dimensional video inputs online?

- Incremental Slow Feature Analysis (IncSFA)

IncSFA for Open-ended RL?

 IncSFA gradually forgets previously learned representations.

IncSFA for Open-ended RL?

 IncSFA gradually forgets previously learned representations.



Curiosity Driven Modular Incremental Slow Feature Analysis



Curiosity Driven Modular Incremental Slow Feature Analysis







Updating Abstraction





Updating Abstraction

Updating Input-Stream Selection Policy






Agent 1. An internal RL agent State Action Noise

- 1. An internal RL agent
- 2. An adaptive IncSFA-ROC abstraction module



- 1. An internal RL agent
- 2. An adaptive IncSFA-ROC abstraction module
- 3. A gating system





Internal Environment



- Internal Environment
- Discrete states (S^{int}) : $\{s_1^{int}, ..., s_n^{int}\}$ equal to the number of input streams



- Internal Environment
- Discrete states (S^{int}) : $\{s_1^{int}, ..., s_n^{int}\}$ equal to the number of input streams
- $\Box \quad \text{At each state } s_i^{\text{int}} \text{ the agent} \\ \text{observes } \tau \text{-samples from } x_i.$



- Internal Environment
- Discrete states (S^{int}): $\{s_1^{\text{int}}, \ldots, s_n^{\text{int}}\}$ equal to the number of input streams
- $\Box \quad \text{At each state } s_i^{\text{int}} \text{ the agent} \\ \text{observes } \tau \text{-samples from } x_i.$

 \Box Actions (A^{int}) : {stay, switch}



Internal Environment

- Discrete states (S^{int}): $\{s_1^{\text{int}}, \ldots, s_n^{\text{int}}\}$ equal to the number of input streams
- $\Box \quad \text{At each state } s_i^{\text{int}} \text{ the agent} \\ \text{observes } \tau \text{-samples from } x_i.$

 \Box Actions (A^{int}) : {stay, switch}



(2) Curious Dr. MISFA's Abstraction Estimator

□ Coupled IncSFA-ROC algorithm

(2) Curious Dr. MISFA's Abstraction Estimator

- □ Coupled IncSFA-ROC algorithm
- □ IncSFA learns slow feature abstractions from input.

(2) Curious Dr. MISFA's Abstraction Estimator

- □ Coupled IncSFA-ROC algorithm
- □ IncSFA learns slow feature abstractions from input.
- □ Robust Online Clustering (ROC; Guedalia et al.) is a clustering algorithm that learns a discrete mapping between the slow feature outputs w.r.t. to the *subjective* state space (S^{\dagger}).

— Online agglomerative clustering algorithm

— Cluster centres merge if they are similar to each other

Subjective State Space (S^{\dagger})

 Proprioceptive states and discrete values of previously encoded abstraction outputs.



Subjective State Space (S^{\dagger})

- Proprioceptive states and discrete values of previously encoded abstraction outputs.
- Multiple nodes of ROC algorithm are used to map slow-feature outputs to the subjective states.



Subjective State Space (S^{\dagger})

- Proprioceptive states and discrete values of previously encoded abstraction outputs.
- Multiple nodes of ROC algorithm are used to map slow-feature outputs to the subjective states.
- Dimensionality of the subjectivestate space increases as the robot acquires more abstractions.





□ For each slow-feature output, the closest ROC node is activated and the corresponding mean-squared estimation error is computed.



- □ For each slow-feature output, the closest ROC node is activated and the corresponding mean-squared estimation error is computed.
- □ A total estimation error is computed as a sum of stored errors corresponding to all the nodes.

Derivative of the total estimation error (ξ) is calculated via backward difference approximation:

$$\xi(t) = \xi(t) - \xi(t-1)$$

Derivative of the total estimation error (ξ) is calculated via backward difference approximation:

• $\xi(t) = \xi(t) - \xi(t-1)$

The agent's reward function is updated at every iteration as follows:

Derivative of the total estimation error (ξ) is calculated via backward difference approximation:

$$\xi$$
 (t) = ξ (t) - ξ (t-1)

The agent's reward function is updated at every iteration as follows:

$$\begin{aligned} R^{\text{int}}(s^{\text{int}}, s^{\text{int}}_{-}, a^{\text{int}}) &= (1 - \eta) \ R^{\text{int}}(s^{\text{int}}, s^{\text{int}}_{-}, a^{\text{int}}) + \eta \sum_{t}^{t + \tau} -\dot{\xi}(t) \end{aligned}$$

where $0 < \eta < 1$ is a discount factor, τ is the number of observed samples (x), $s^{\text{int}}, s^{\text{int}}_{-}) \in \mathcal{S}^{\text{int}}$ and $a^{\text{int}} \in \{\text{stay, switch}\}. \end{aligned}$



 $\Box \pi^{\text{int}}: S^{\text{int}} \to \{stay, switch\}$

- $\Box \ \pi^{\text{int}}: S^{\text{int}} \rightarrow \{stay, switch\}$
- □ Learned through Least-Squares Temporal Difference (LSTD) approach using the current estimate of the reward function.

- $\Box \pi^{\text{int}}: S^{\text{int}} \to \{stay, switch\}$
- □ Learned through Least-Squares Temporal Difference (LSTD) approach using the current estimate of the reward function.
- Agent explores using a decaying epsilon-greedy strategy over the inputstream selection policy.

(c) Gating System



□ Curiosity Function: Quantifies learning difficulty of the abstraction estimator $(\Omega : X \rightarrow [0,1])$

- □ Curiosity Function: Quantifies learning difficulty of the abstraction estimator $(\Omega : X \rightarrow [0,1])$
- \square Ω induces a total ordering among the input streams.

- □ Curiosity Function: Quantifies learning difficulty of the abstraction estimator $(\Omega : X \rightarrow [0,1])$
- \square Ω induces a total ordering among the input streams.
- \square Easier to encode input streams have lower Ω values.

- Curiosity Function: Quantifies learning difficulty of the abstraction estimator $(\Omega : X \rightarrow [0,1])$
- \square Ω induces a total ordering among the input streams.
- \square Easier to encode input streams have lower Ω values.

Definition 1 We define the curiosity-function of the IncSFA algorithm for an input stream $\mathbf{x}(t) \in \mathbb{R}^{I}$ as

$$\Omega(\mathbf{x}) = \left[1 - \frac{\eta^{mca}(\lambda_{K-1} - \lambda_K)}{1 - \eta^{mca} - \eta^{mca}\lambda_K}\right]$$
(23)

where λ_K and λ_{K-1} denote the smallest two eigenvalues of $E[\dot{\mathbf{z}}(t)\dot{\mathbf{z}}(t)^T]$, $\mathbf{z}(t) \in \mathbb{R}^K$ is the whitened output of $\mathbf{x}(t)$.









X1

X2

Abstraction Learning to Skill Acquisition



Abstraction Learning to Skill Acquisition



Skill Acquisition using Curious Dr. MISFA (Kompella et al. 2014)

Task independent curiosity-driven skill acquisition algorithm
Skill Acquisition using Curious Dr. MISFA (Kompella et al. 2014)

- □ Task independent curiosity-driven skill acquisition algorithm
- □ Combines Curious Dr. MISFA with the options Framework
 - Curious Dr. MISFA Extracts slow feature abstractions as quickly as possible
 - Options Builds the abstractions into skills simultaneously

Options Framework

Option: Temporally extended courses of actions.



Options Framework

- Option: Temporally extended courses of actions.
- □ Formally defined as a tuple $O = \langle I, \beta, \pi \rangle$:
 - $I \subseteq S$ (Initiation set),

 $\beta: S \rightarrow [0,1]$ (Termination cond.),

 $\pi: I \times A \rightarrow [0,1]$ (Option policy)



Options Framework

- Option: Temporally extended courses of actions.
- □ Formally defined as a tuple $O = \langle I, \beta, \pi \rangle$:
 - $I \subseteq S$ (Initiation set),

 $\beta: S \rightarrow [0,1]$ (Termination cond.),

 π : I × *A* \rightarrow [0,1] (Option policy)

 $\Box \quad \text{Executing option} \rightarrow \text{High-dimensional observations}$





Random exploratory → Deterministic and increasingly complex behaviors target behaviors

(Exploratory Options)



Random exploratory→Deterministic and increasingly complex
target behaviors

(Exploratory Options)



Random exploratory→Deterministic and increasingly complex
target behaviors

(Exploratory Options)



Random exploratory → Deterministic and increasingly complex behaviors target behaviors

(Exploratory Options)

Exploratory Options (O^e)

 $\Box \quad O^e = \langle I^e, \beta^e, \pi^e \rangle$

 $\Box \quad \mathbf{I}^{\mathbf{e}} \subseteq S^{\dagger}$

- π^e: Exploratory-option's stochastic policy is a random walk in option's state-space
- $\square \quad \beta^{e}: \text{The option terminates} \\ \text{after } \boldsymbol{\tau} \text{ time-steps since its} \\ \text{execution.} \\ \end{bmatrix}$



Target Options (Acquired Skills)



Target Options (Acquired Skills)

 $\Box \quad O^{\mathcal{L}} = \langle I^{\mathcal{L}}, \beta^{\mathcal{L}}, \phi, \pi^{\mathcal{L}} \rangle$

Target Options (Acquired Skills)

- $\Box \quad O^{\mathcal{L}} = \left< I^{\mathcal{L}}, \, \beta^{\mathcal{L}}, \, \phi \,, \, \pi^{\mathcal{L}} \right>$
- Target option's policy maximises the variation of observed featureoutputs.



Learning a Target Option (0²)

 \Box Say Curious Dr. MISFA learns an abstraction ϕ_i corresponding to x_j

- □ Say Curious Dr. MISFA learns an abstraction ϕ_i corresponding to x_j
- $\square \quad \text{Initiation set: } \mathbf{I}_{i}^{\mathcal{L}} = (\mathbf{I}_{i}^{e} \times S^{\phi_{i}})$

- Say Curious Dr. MISFA learns an abstraction ϕ_i corresponding to x_j
- $\square \quad \text{Initiation set: } \mathbf{I}_{i}^{\mathcal{L}} = (\mathbf{I}_{i}^{e} \times S^{\phi_{i}})$
- Target option policy ($\pi_i^{\mathcal{L}} : \mathbf{I}_i^{\mathcal{L}} \to A$) : Developed through LSPI using:
 - Estimated option's transition model from samples generated by π_{j}^{e} - Estimated options' reward model using diff. of subseq. abstraction activations:

- Say Curious Dr. MISFA learns an abstraction ϕ_i corresponding to x_j
- $\Box \quad \text{Initiation set: } \mathbf{I}_{i}^{\mathcal{L}} = (\mathbf{I}_{i}^{e} \times S^{\phi_{i}})$
- Target option policy ($\pi_i^{\mathcal{L}} : \mathbf{I}_i^{\mathcal{L}} \to A$) : Developed through LSPI using:
 - Estimated option's transition model from samples generated by π_{j}^{e} - Estimated options' reward model using diff. of subseq. abstraction activations:

 $R^{O_i^{\mathcal{L}}}(s^{\dagger}, a, s_{-}^{\dagger}) = (1 - \alpha) R^{O_i^{\mathcal{L}}}(s^{\dagger}, a, s_{-}^{\dagger}) + \alpha \|\mathbf{y}_{\mathbf{i}}(t) - \mathbf{y}_{\mathbf{i}}(t - 1)\|$ where $\mathbf{y}_{\mathbf{i}}(t) = \phi_i \left(\mathcal{U}(\mathcal{P}(s_{-}, \pi_j^e(s_{-}^{\dagger}))) \right)$ and $\mathbf{y}_{\mathbf{i}}(t - 1) = \phi_i \left(\mathcal{U}(\mathcal{P}(s, \pi_j^e(s^{\dagger}))) \right)$

- Say Curious Dr. MISFA learns an abstraction ϕ_i corresponding to x_j
- $\Box \quad \text{Initiation set: } \mathbf{I}_{i}^{\mathcal{L}} = (\mathbf{I}_{i}^{e} \times S^{\phi_{i}})$
- □ Target option policy ($\pi_i^{\mathcal{L}}$: $\mathbf{I}_i^{\mathcal{L}} \to A$) : Developed through LSPI using:
 - Estimated option's transition model from samples generated by π_{j}^{e} - Estimated options' reward model using diff. of subseq. abstraction activations:

 $R^{O_i^{\mathcal{L}}}(s^{\dagger}, a, s_{-}^{\dagger}) = (1 - \alpha) R^{O_i^{\mathcal{L}}}(s^{\dagger}, a, s_{-}^{\dagger}) + \alpha \|\mathbf{y}_{\mathbf{i}}(t) - \mathbf{y}_{\mathbf{i}}(t - 1)\|$ where $\mathbf{y}_{\mathbf{i}}(t) = \phi_i \left(\mathcal{U}(\mathcal{P}(s_{-}, \pi_j^e(s_{-}^{\dagger}))) \right)$ and $\mathbf{y}_{\mathbf{i}}(t - 1) = \phi_i \left(\mathcal{U}(\mathcal{P}(s, \pi_j^e(s^{\dagger}))) \right)$

Termination Condition: Option terminates whenever the agent observes the maximum of estimated reward model.

Summary of Policies Involved

 \square π^{int} : Internal Policy that is learned to determine which exploratory option to execute.

Summary of Policies Involved

- \square π^{int} : Internal Policy that is learned to determine which exploratory option to execute.
- $\square \pi^{e}$: Exploratory-option's stochastic policy that is used to observe highdimensional observations to update slow feature abstraction.

Summary of Policies Involved

- $\square \pi^{\text{int}}$: Internal Policy that is learned to determine which exploratory option to execute.
- $\square \pi^{e}$: Exploratory-option's stochastic policy that is used to observe highdimensional observations to update slow feature abstraction.
- \square $\pi^{\mathcal{L}}$: Target-option's deterministic policy that is learned to maximize variation in the slow-feature abstraction output.

Experimental Results



iCub's Left-Camera Image

iCub's Right-Camera Image





iCub's state space dimensionality

Vision: 640 x 480 pixels Joints: 51 motors Total = 640 x 480 x 51 = (15,667,200) states variables!

Our Focus:





Low-dimensional joint state space embedding learned apriori



Experiment 1 iCub Learns to Topple a Cup

https://www.youtube.com/watch?v=OTqdXbTEZpE

https://www.youtube.com/watch?v=OTqdXbTEZpE



Experiment 2

iCub Learns to Grasp a Cup

https://www.youtube.com/watch?v=OTqdXbTEZpE

https://www.youtube.com/watch?v=OTqdXbTEZpE



le Input Observation Streams (X) (left came

- □ Simple Domains:
 - Bakker & Schmidhuber [2004] (Hassle; Hierarchical RL)
 - □ Stout & Barto [2010] (Competence-based IM)
 - □ Pape et al. [2012] (tactile skills on a biomimetic finger)

Attempts to find skills using feature-abstractions in the domain of Humanoid robots:

 \square Hart [2009] - (DEDS + task-specific IM)

- Attempts to find skills using feature-abstractions in the domain of Humanoid robots:
 - \Box Hart [2009] (DEDS + task-specific IM)
 - Konidaris [2009-2011] (Options + ART tags + pre-existing abstraction library), uBot.

- Attempts to find skills using feature-abstractions in the domain of Humanoid robots:
 - \Box Hart [2009] (DEDS + task-specific IM)
 - Konidaris [2009-2011] (Options + ART tags + pre-existing abstraction library), uBot.
 - □ Mugan & Kuipers [2012] (QLAP, assumes low-level object tracking models)

- Attempts to find skills using feature-abstractions in the domain of Humanoid robots:
 - \Box Hart [2009] (DEDS + task-specific IM)
 - Konidaris [2009-2011] (Options + ART tags + pre-existing abstraction library), uBot.
 - □ Mugan & Kuipers [2012] (QLAP, assumes low-level object tracking models)
 - Barnes & Oudeyer [2013] (SAGG-RIAC, Competence progress + heuristics + subgoals, task-dependent)

- Attempts to find skills using feature-abstractions in the domain of Humanoid robots:
 - \Box Hart [2009] (DEDS + task-specific IM)
 - Konidaris [2009-2011] (Options + ART tags + pre-existing abstraction library), uBot.
 - □ Mugan & Kuipers [2012] (QLAP, assumes low-level object tracking models)
 - □ Barnes & Oudeyer [2013] (SAGG-RIAC, Competence progress + heuristics + subgoals, task-dependent)
 - □ Ngo et al. [2012-2013] (selective sampling + Curiosity, Overhead camera + pre-designed sensori-motor abstractions)
Take-Home Message



Take-Home Message

□ Curious Dr. MISFA is an online active modular incremental SFA.

Take-Home Message

- □ Curious Dr. MISFA is an online active modular incremental SFA.
- Skill acquisition using Curious Dr. MISFA is a task-independent method that demonstrates intrinsically motivated skill acquisition from raw-pixel streams on a real humanoid robot.

Acknowledgement

 This tutorial is funded through SNF grant #138219 (Theory and Practice of Reinforcement Learning II)

I personally thank:

Dr. Jan Koutnik (<u>http://www.idsia.ch/~koutnik/</u>), Dr. Faustino Gomez (<u>http://www.idsia.ch/~tino/</u>) and Dr. Matthew Luciw (<u>http://www.idsia.ch/~luciw/</u>) for helping me prepare for this tutorial.

Disclaimer

- Some of the artwork used in this presentation was downloaded from Google Images.
- Presenter does not wish to use them for any commercial purposes.



Thank you